Supplementary Material of Linear Label Ranking with Bounded Noise

Anonymous Author(s) Affiliation Address email

Abstract

Label Ranking (LR) is the supervised task of learning a sorting function that maps 1 feature vectors $x \in \mathbb{R}^d$ to rankings $\sigma(x) \in \mathbb{S}_k$ over a finite set of k labels. We 2 focus on the fundamental case of learning linear sorting functions (LSFs) under 3 Gaussian marginals: x is sampled from the d-dimensional standard normal and 4 the ground truth ranking $\sigma^{\star}(x)$ is the ordering induced by sorting the coordinates 5 of the vector $W^{\star}x$, where $W^{\star} \in \mathbb{R}^{k \times d}$ is unknown. We consider learning 6 LSFs in the presence of bounded noise: assuming that a noiseless example is of 7 the form $(x, \sigma^*(x))$, we observe (x, π) , where for any pair of elements $i \neq j$, 8 the probability that the order of i, j is different in π than in $\sigma^{\star}(x)$ is at most 9 $\eta < 1/2$. We design efficient non-proper and proper learning algorithms that 10 learn hypotheses within normalized Kendall's Tau distance ϵ from the ground truth 11 with $N = O(d \log(k)/\epsilon)$ labeled examples and runtime poly(N, k). For the more 12 challenging top-r disagreement loss, we give an efficient proper learning algorithm 13 that achieves ϵ top-r disagreement with the ground truth with $N = O(dkr/\epsilon)$ 14 samples and poly(N) runtime. 15

16 1 Introduction

17 1.1 Background and Motivation

Label Ranking (LR) is the problem of learning a hypothesis that maps features to rankings over a 18 finite set of labels. Given a feature vector $x \in \mathbb{R}^d$, a sorting function $\sigma(\cdot)$ maps it to a ranking of k 19 alternatives, i.e., $\sigma(x)$ is an element of the symmetric group with k elements, \mathbb{S}_k . Assuming access 20 to a training dataset of features labeled with their corresponding rankings, i.e., pairs of the form 21 $(x, \pi) \in \mathbb{R}^d \times \mathbb{S}_k$, the goal of the learner is to find a sorting function h(x) that generalizes well over 22 a fresh sample. LR has received significant attention over the years [DSM03, SS07, HFCB08, CH08, 23 FHMB08] due to the large number of applications. For example, ad targeting [DGR⁺14] is an LR 24 instance where for each user we want to use their feature vector to predict a ranking over ad categories 25 and present them with the most relevant. The practical significance of LR has lead to the development 26 of many techniques based on probabilistic models and instance-based methods [CH08, CDH10], 27 [GDV12, ZLGQ14], decision trees [CHH09], entropy-based ranking trees [RdSRSK15], bagging 28 [AGM17], and random forests [dSSKC17, ZQ18]. However, almost all of these works come without 29 provable guarantees and/or fail to learn in the presence of noise in the observed rankings. 30

Linear Sorting Functions (LSFs). In this work, we focus on the fundamental concept class of Linear Sorting functions [HPRZ03]. A linear sorting function parameterized by a matrix $W \in \mathbb{R}^{k \times d}$ with k rows W_1, \ldots, W_k takes a feature $x \in \mathbb{R}^d$, maps it to $Wx = (W_1 \cdot x, \ldots, W_k \cdot x) \in \mathbb{R}^k$ and then outputs an ordering (i_1, \ldots, i_k) of the k alternatives such that $W_{i_1} \cdot x \ge W_{i_2} \cdot x \ge \ldots \ge W_{i_k} \cdot x$.

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

In other words, a linear sorting function ranks the k alternatives (corresponding to rows of W) with respect to how well they correlate with the feature x. We denote a linear sorting function with parameter $W \in \mathbb{R}^{k \times d}$ by $\sigma_W(x) \triangleq \operatorname{argsort}(Wx)$ where $\operatorname{argsort} : \mathbb{R}^k \to \mathbb{S}_k$ takes as input a vector $(v_1, \ldots, v_k) \in \mathbb{R}^k$, sorts it in decreasing order to obtain $v_{i_1} \ge v_{i_2} \ge \ldots \ge v_{i_k}$ and returns the ordering (i_1, \ldots, i_k) .

40 **Noisy Ranking Distributions.** Learning LSFs in the noiseless setting can be done efficiently by 41 using linear programming. However, the common assumption both in theoretical and in applied 42 works is that the observed rankings are noisy in the sense that they do not always correspond to 43 the ground-truth ranking. We assume that the probability that the order of two elements i, j in the 44 observed ranking π is different than their order in the ground-truth ranking σ^* is at most $\eta < 1/2$.

45 **Definition 1** (Noisy Ranking Distribution). Fix $\eta \in [0, 1/2)$. An η -noisy ranking distribution 46 $\mathcal{M}(\sigma^*)$ with ground-truth ranking $\sigma^* \in \mathbb{S}_k$ is a probability measure over \mathbb{S}_k that, for any $i, j \in [k]$, 47 with $i \neq j$, satisfies $\mathbf{Pr}_{\pi \sim \mathcal{M}(\sigma^*)}[i \prec_{\pi} j \mid i \succ_{\sigma^*} j] \leq \eta$.¹

Note that, when $\eta = 0$, we always observe the ground-truth permutation and, in the case of $\eta =$ 48 1/2, we may observe a uniformly random permutation. We remark that most natural ranking 49 distributions satisfy this bounded noise property, e.g., (i) the Mallows model, which is probably 50 the most fundamental ranking distribution (see, e.g., [BM09, LB11, CPS13, ABSV14, BFFSZ19, 51 FKS21, DOS18, LM18, MW20, LM21] for a small sample of this line of research) and (ii) the 52 Bradley-Terry-Mallows model [Mal57], which corresponds to the ranking distribution analogue of 53 the Bradley-Terry-Luce model [BT52, Luc12] (the most studied pairwise comparisons model; see, 54 e.g., [Hun04, NOS17, APA18] and the references therein). For more details, see Supp. Material E. 55

⁵⁶ We consider the fundamental setting where the feature vector $x \in \mathbb{R}^d$ is generated by a standard ⁵⁷ normal distribution and the ground-truth ranking for each sample x is given by the LSF $\sigma_{W^*}(x)$ for ⁵⁸ some unknown parameter matrix $W^* \in \mathbb{R}^{k \times d}$. For a fixed x, the ranking that we observe comes ⁵⁹ from an η -noisy ranking distribution with ground-truth ranking $\sigma_{W^*}(x)$.

60 **Definition 2** (Noisy Linear Label Ranking Distribution). Fix $\eta \in [0, 1/2)$ and some ground-truth 61 parameter matrix $\mathbf{W}^* \in \mathbb{R}^{d \times k}$. We assume that the η -noisy linear label ranking distribution \mathcal{D} over 62 $\mathbb{R}^d \times \mathbb{S}_k$ satisfies the following:

1. The *x*-marginal of \mathcal{D} is the *d*-dimensional standard normal distribution.

64 2. For any $(x, \pi) \sim D$, the distribution of π conditional on x is an η -noisy ranking distribution 65 with ground-truth ranking $\sigma_{W^*}(x)$.

At first sight, the assumption that the underlying x-marginal is the standard normal may look too 66 strong. However, for k = 2, Definition 2 captures the problem of learning linear threshold functions 67 with Massart noise. Without assumptions for the x-marginal, it is known [DGT19, CKMY20, DK20, 68 NT22] that optimal learning of halfspaces under Massart noise requires super-polynomial time (in 69 the Statistical Query model of [Kea98]). On the other hand, a lot of recent works [BZ17, MV19, 70 DKTZ20, ZSA20, ZL21] have obtained efficient algorithms for learning Massart halfspaces under 71 72 Gaussian marginals. The goal of this work is to provide efficient algorithms for the more general 73 problem of learning LSFs with bounded noise under Gaussian marginals.

74 **1.2 Our Results**

The main contributions of this paper are the first efficient algorithms for learning LSFs with bounded noise with respect to Kendall's Tau distance and top-r disagreement loss.

⁷⁷ Learning in Kendall's Tau Distance. The most standard metric in rankings [SSBD14] is Kendall's ⁷⁸ Tau (KT) distance which, for two rankings $\pi, \tau \in S_k$, measures the fraction of pairs (i, j) on ⁷⁹ which they disagree. That is, $\Delta_{\text{KT}}(\pi, \tau) = \sum_{i \prec \pi j} \mathbf{1}\{i \succ_{\tau} j\}/{\binom{k}{2}}$. Our first result is an efficient ⁸⁰ learning algorithm that, given samples from an η -noisy linear label ranking distribution \mathcal{D} , computes ⁸¹ a parameter matrix W that ranks the alternatives almost optimally with respect to the KT distance ⁸² from the ground-truth ranking $\sigma_{W^*}(\cdot)$.

¹We use $i \succ_{\pi} j$ (resp. $i \prec_{\pi} j$) to denote that the element i is ranked higher (resp. lower) than j according to the ranking π .

- **Theorem 1** (Learning LSFs in KT Distance). Fix $\eta \in [0, 1/2)$ and $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy 83
- linear label ranking distribution satisfying the assumptions of Definition 2 with ground-truth LSF 84
- 85
- $\sigma_{\mathbf{W}^{\star}}(\cdot)$. There exists an algorithm that draws $N = \tilde{O}\left(\frac{d}{\epsilon(1-2\eta)^6}\log(k/\delta)\right)$ samples from \mathcal{D} , runs in sample-polynomial time, and computes a matrix $\mathbf{W} \in \mathbb{R}^{k \times d}$ such that, with probability at least 86

$$\mathop{\mathbf{E}}_{\boldsymbol{x}\sim\mathcal{N}_d}[\Delta_{\mathrm{KT}}(\sigma_{\boldsymbol{W}}(\boldsymbol{x}),\sigma_{\boldsymbol{W}^{\star}}(\boldsymbol{x}))] \leq \epsilon$$

Theorem 1 gives the first efficient algorithm with provable guarantees for the supervised problem of 88 learning noisy linear rankings. We remark that the sample complexity of our learning algorithm is 89 qualitatively optimal (up to logarithmic factors) since, for k = 2, our problem subsumes learning 90 a linear classifier with Massart noise ² for which $\Omega(d/\epsilon)$ are known to be information theoretically 91 necessary [MN06]. Moreover, our learning algorithm is *proper* in the sense that it computes a 92 linear sorting function $\sigma_{\mathbf{W}}(\cdot)$. As opposed to improper learners (see also Section 1.3), a proper 93 learning algorithm gives us a compact representation (storing W requires O(kd) memory) of the 94 sorting function that allows us to efficiently compute (with runtime $O(kd + k \log k)$) the ranking 95 corresponding to a fresh datapoint $x \in \mathbb{R}^d$. 96

Learning in top-r Disagreement. We next present our learning algorithm for the top-r metric 97 formally defined as $\Delta_{top-r}(\pi, \tau) = \mathbf{1}\{\pi_{1..r} \neq \tau_{1..r}\}$, where by $\pi_{1..r}$ we denote the ordering on the first *r* elements of the permutation π . The top-*r* metric is a disagreement metric in the sense that it 98 99 takes binary values and for r = 1 captures the standard (multiclass) top-1 classification loss. We 100 remark that, in contrast with the top-r classification loss, which only requires the predicted label to 101 be in the top-r predictions of the model, the top-r ranking metric that we consider here requires that 102 the model puts *the same elements in the same order* as the ground truth in the top-r positions. The 103 top-r ranking is well-motivated as, for example, in ad targeting (discussed in Section 1.1) we want to 104 105 be accurate on the top-r ad categories for a user so that we can diversify the content that they receive. **Theorem 2** (Learning LSFs in top-r Disagreement). Fix $\eta \in [0, 1/2)$, $r \in [k]$ and $\epsilon, \delta \in (0, 1)$. 106

Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2 107 with ground-truth LSF $\sigma_{\mathbf{W}^{\star}}(\cdot)$. There exists an algorithm that draws $N = \widetilde{O}\left(\frac{drk}{\epsilon(1-2\eta)^6}\log(1/\delta)\right)$ samples from \mathcal{D} , runs in sample-polynomial time and computes a matrix $\mathbf{W} \in \mathbb{R}^{k \times d}$ such that, with 108

109 probability at least $1 - \delta$. 110

$$\mathop{\mathbf{E}}_{\boldsymbol{x}\sim\mathcal{N}_d}[\Delta_{\mathrm{top}-r}(\sigma_{\boldsymbol{W}}(\boldsymbol{x}),\sigma_{\boldsymbol{W}^{\star}}(\boldsymbol{x}))] \leq \epsilon.$$

As a direct corollary of our result, we obtain a proper algorithm for learning the top-1 element 111 with respect to the standard 0-1 loss that uses O(kd) samples. In fact, for small values of r, i.e., 112 r = O(1), our sample complexity is essentially tight. It is known that $\Theta(kd)$ samples are information 113 theoretically necessary [Nat89] for top-1 classification. ³ For the case r = k, i.e., when we want to 114 learn the whole ranking with respect to the 0-1 loss, our sample complexity is $O(k^2 d)$. However, 115 using arguments similar to [DSBDSS11], one can show that in fact O(dk) ranking samples are 116 sufficient in order to learn the whole ranking with respect to the 0-1 loss. In this case, it is unclear 117 whether a better sample complexity can be achieved with an efficient algorithm and we leave this as 118 an interesting open question for future work. 119

1.3 Our Techniques 120

 $1-\delta$

87

Learning in Kendall's Tau distance. Our proper learning algorithm consists of two steps: an 121 improper learning algorithm that decomposes the ranking problem to $O(k^2)$ binary linear classifica-122 tion problems and a convex (second order conic) program that "compresses" the k^2 linear classifiers 123 to obtain a $k \times d$ matrix W. Our improper learning algorithm splits the ranking learning problem 124 into $O(k^2)$ binary, d-dimensional linear classification problems with Massart noise. In particular, 125 for every pair of elements $i, j \in [k]$, each binary classification task asks whether element i is 126

²Notice that in this case Kendall's Tau distance is simply the standard 0-1 binary loss.

³Strictly speaking, those lower bounds do not directly apply in our setting because our labels are whole rankings instead of just the top classes but, in the Supp. Material D, we show that we can adapt the lower bound technique of [DSBDSS11] to obtain the same sample complexity lower bound for our ranking setting.

ranked higher than element j in the ground-truth permutation $\sigma_{W^*}(x)$. As we already discussed, 127 we have that, under the Gaussian distribution, there exist efficient Massart learning algorithms 128 [BZ17, MV19, DKTZ20, ZSA20, ZL21] that can recover linear classifiers sgn $(v_{ij} \cdot x)$ that correctly 129 order the pair i, j for all x apart from a region of $O(\epsilon)$ -Gaussian mass. However, we still need 130 to aggregate the results of the *approximate* binary classifiers in order to obtain a ranking of the k131 alternatives for each x. We first show that we can design a "voting scheme" that combines the results 132 of the binary classifiers using an efficient constant factor approximation algorithm for the Minimum 133 Feedback Arc Set (MFAS) problem [ACN08]. This gives us an efficient but improper algorithm for 134 learning LSFs in Kendall's Tau distance. In order to obtain a proper learning algorithm, we further 135 "compress" the $O(k^2)$ approximate linear classifiers with normal vectors v_{ij} and obtain a matrix 136 $W \in \mathbb{R}^{k \times d}$ with the property that the difference of every two rows $W_i - W_j$ is $O(\epsilon)$ -close to the vector v_{ij} . More precisely, we show that, given the linear classifiers $v_{ij} \in \mathbb{R}^d$, we can efficiently compute a matrix $W \in \mathbb{R}^{k \times d}$ such that the following angle distance with W^* is small: 137 138 139

$$d_{\text{angle}}(\boldsymbol{W}, \boldsymbol{W}^{\star}) \triangleq \max_{i,j} \theta(\boldsymbol{W}_i - \boldsymbol{W}_j, \boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}) \le O(\epsilon) \,. \tag{1}$$

It is not hard to show that, as long as the above angle metric is at most $O(\epsilon)$, then (in expectation over the standard Gaussian) Kendall's Tau distance between the LSFs is also $O(\epsilon)$. A key technical difficulty that we face in this reduction is bounding the "condition number" of the convex (second order conic) program that finds the matrix W given the vectors v_{ij} , see Claim 2. Finally, we remark that the proper learning algorithm of Theorem 1 results in a compact and efficient sorting function that requires: (i) storing O(k) weight vectors as opposed to the initial $O(k^2)$ vectors of the improper learner; and (ii) evaluating k inner products with x to find its ranking (instead of $O(k^2)$).

Learning in top-r Disagreement. We next turn our attention to the more challenging top-r ranking 147 disagreement metric. In particular, suppose that we are interested in recovering only the top element 148 of the ranking. One approach would be to directly use the improper learning algorithm for this 149 task and ask for KT distance of order roughly ϵ/k^2 . The resulting hypothesis would produce good 150 predictions for the top element but the required sample complexity would be $O(dk^2)$. While it seems 151 that training $O(k^2)$ d-dimensional binary classifiers inherently requires $O(dk^2)$ samples, we show 152 that, using the proper KT distance learning algorithm of Theorem 1, we can also obtain improved 153 sample complexity results for the top-r metric. Our main technical contribution here is a novel 154 estimate of the top-r disagreement in terms of the angle metric. In general, one can show that the top-r disagreement is at most $O(k^2) d_{\text{angle}}(\boldsymbol{W}, \boldsymbol{W}^*)$. We significantly sharpen this estimate by 155 156 showing the following lemma. 157

Lemma 1 (Top-*r* Disagreement via Parameter Distance). Consider two matrices $W, W^* \in \mathbb{R}^{k \times d}$ and let \mathcal{N}_d be the standard Gaussian in d dimensions. We have that

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[\sigma_{1..r}(\boldsymbol{W}\boldsymbol{x})\neq\sigma_{1..r}(\boldsymbol{W}^{\star}\boldsymbol{x})]\leq \widetilde{O}(kr)\;d_{\mathrm{angle}}(\boldsymbol{W},\boldsymbol{W}^{\star})\,.$$

We remark that Lemma 1 is a general geometric tool that we believe will be useful in other distributionspecific multiclass learning settings. The proof of Lemma 1 mainly relies on geometric Gaussian surface area computations that we believe are of independent interest. For the details, we refer the reader to Section 4. An interesting question with a convex-geometric flavor is whether the sharp bound of Lemma 1 also holds under the more general class of isotropic log-concave distributions.

165 1.4 Related Work

Robust Supervised Learning. We start with a summary of prior work on PAC learning with Massart 166 noise. The Massart noise model was formally defined in [MN06] but similar variants had been defined 167 by Vapnik, Sloan and Rivest [Vap06, Slo88, Slo92, RS94, Slo96]. This model is a strict extension 168 of the Random Classification Noise (RCN) model [AL88], where the label noise is uniform, i.e., 169 context-independent and is a special case of the agnostic model [Hau18, KSS94], where the label 170 noise is fully adversarial and computational barriers are known to exist [GR09, FGKP06, Dan16, 171 DKZ20, GGK20, DKPZ21, HSSVG22]. Our work partially builds upon on the algorithmic task of 172 PAC learning halfspaces with Massart noise [BH20]. In the distribution-independent setting, known 173 efficient algorithms [DGT19, CKMY20, DKT21] achieve error $\eta + \epsilon$ and the works of [DK20, NT22] 174 indicate that this error bound is the best possible in the Statistical Query model [Kea98]. This lower 175

bound motivates the study of the distribution-specific setting (which is also the case of our work). 176 There is an extensive line of work in this direction: [ABHU15, ABHZ16, YZ17, ZLC17, BZ17, 177 MV19, DKTZ20, ZSA20, ZL21] with the currently best algorithms succeeding for all $\eta < 1/2$ 178 with a sample and computational complexity $poly(d, 1/\epsilon, 1/(1-2\eta))$ under a class of distributions 179 including isotropic log-concave distributions. For details, see [DKK⁺21]. In this work we focus on 180 Gaussian marginals but some of our results extend to larger distribution classes. 181 Label Ranking. Our work lies in the area of Label Ranking, which has received significant attention 182 over the years [SS07, HFCB08, CH08, HPRZ03, FHMB08, DSM03]. There are multiple approaches 183 for tackling this problem (see $[VG10], [ZLY^+14]$). Some of them are based on probabilistic models 184 [CH08, CDH10, GDV12, ZLGQ14] or may be tree based, such as decision trees [CHH09], entropy 185 based ranking trees and forests [RdSRSK15, dSSKC17], bagging techniques [AGM17] and random 186 forests [ZQ18]. There are also works focusing on supervised clustering [GDGV13]. Finally, [CH08, 187 CDH10, CHH09] adopt an instance-based approaches using nearest neighbors approaches. The above 188 results are industrial. From a theoretical perspective, LR has been mainly studied from a statistical 189 learning theory framework [CV20, CKS18, KGB18, KCS17]. [FKP21] provide some computational 190 guarantees for the performance of decision trees in the noiseless case and some experimental results 191 on the robustness of random forests to noise. The setting of $[DGR^{+}14]$ is close to ours but is 192 investigated from an experimental standpoint. We remark that while reducing LR to multiple binary 193 classification tasks has been used in prior literature [HFCB08, CH12, FKP21], standard reductions 194 can not tolerate noise in rankings (nevertheless, from an experimental perspective, e.g., random 195 forests seem robust to noise but lack formal theoretical guarantees). Our reduction crucially relies on 196 the existence of efficient learning algorithms for binary linear classification with Massart noise. 197

198 2 Notation and Preliminaries

General Notation. We use $O(\cdot)$ to omit poly-logarithmic factors. A learning algorithm has samplepolynomial runtime if it runs in time polynomial in the size of the description of the input training set. We denote vectors by boldface \boldsymbol{x} (with elements x_i) and matrices with \boldsymbol{W} , where we let $\boldsymbol{W}_i \in \mathbb{R}^d$ denote the *i*-th row of $\boldsymbol{W} \in \mathbb{R}^{k \times d}$ and W_{ij} its elements. We denote $\boldsymbol{a} \cdot \boldsymbol{b}$ the inner product of two vectors and $\theta(\boldsymbol{a}, \boldsymbol{b})$ their angle. Let \mathcal{N}_d denote the *d*-dimensional standard normal and $\Gamma(\cdot)$ the Gaussian surface area.

Rankings. We let $\operatorname{argsort}_{i \in [k]} v$ denote the ranking of [k] in decreasing order according to the values of v. For a ranking π , we let $\pi(i)$ denote the position of the *i*-th element. If $\pi = \pi(x)$, we may also write $\pi(x)(i)$ to denote the position of *i*. We often refer to the elements of a ranking as *alternatives*. For a ranking σ , we let $\sigma_{1..r}$ denote the top-*r* part of σ . When $\sigma = \sigma(x)$, we may also write $\sigma_{1..r}(x)$ and $\sigma_{\ell}(x)$ will be the alternative at the ℓ -th position. We let Δ_{KT} denote the (normalized) KT distance, i.e., $\Delta_{\mathrm{KT}}(\pi, \tau) = \sum_{i \prec \pi j} \mathbf{1}\{i \succ_{\tau} j\}/{\binom{k}{2}}$ for $\pi, \tau \in \mathbb{S}_k$.

3 Learning in KT distance: Theorem 1

In this section, we present the main tools required to obtain our proper learning algorithm of Theorem 1. Our proper algorithm adopts a two-step approach: it first invokes an efficient *improper* algorithm which, instead of a linear sorting function (i.e., a matrix $W \in \mathbb{R}^{k \times d}$), outputs a list of $O(k^2)$ linear classifiers. We then design a novel convex program in order to find the matrix Wsatisfying the guarantees of Theorem 1. Let us begin with the improper learner for LSFs with bounded noise with respect to the KT distance, whose description can be found in Algorithm 1.

218 3.1 Improper Learning Algorithm

Let us assume that the target function is $\sigma^{\star}(x) = \sigma_{W^{\star}}(x) = \operatorname{argsort}(W^{\star}x)$ for some $W^{\star} \in \mathbb{R}^{k \times d}$.

Step 1: Binary decomposition and Noise Structure. For each drawn example (x, π) from the η -noisy linear label ranking distribution \mathcal{D} (see Definition 2), we create $\binom{k}{2}$ binary examples (x, y_{ij}) with $y_{ij} = \text{sgn}(\pi(i) - \pi(j))$ for any $1 \le i < j \le k$. We have that

$$\Pr_{(\boldsymbol{x}, \pi) \sim \mathcal{D}} \left[y_{ij} \cdot \operatorname{sgn}((\boldsymbol{W}_i^\star - \boldsymbol{W}_j^\star) \cdot \boldsymbol{x}) < 0 \mid \boldsymbol{x} \right] = \Pr_{\pi \sim \mathcal{M}(\sigma^\star(\boldsymbol{x}))} \left[\pi(i) < \pi(j) \mid \boldsymbol{W}_i^\star \cdot \boldsymbol{x} < \boldsymbol{W}_j^\star \cdot \boldsymbol{x} \right] \,.$$

Algorithm 1 Non-proper Learning Algorithm ImproperLSF

Input: Training set $T = \{(\boldsymbol{x}^t, \pi^t)\}_{t \in [N]}, \epsilon, \delta \in (0, 1), \eta \in [0, 1/2)$ Output: Sorting function $h : \mathbb{R}^d \to \mathbb{S}_k$

 $\begin{array}{ll} \text{For any } 1 \leq i < j \leq k \text{, create } T_{ij} = \{(\boldsymbol{x}^t, \text{sgn}(\pi^t(i) - \pi^t(j)))\} \\ \text{For any } 1 \leq i < j \leq k \text{, compute } \boldsymbol{v}_{ij} = \texttt{MassartLTF}(T_{ij}, \frac{\epsilon}{4}, \frac{\delta}{10k^2}, \eta) \\ \text{Ranking Phase: Given } \boldsymbol{x} \in \mathbb{R}^d \text{:} \\ \text{(a) Construct directed graph } G \text{ with } V(G) = [k] \text{ and edges } e_{i \rightarrow j} \text{ only if } \boldsymbol{v}_{ij} \cdot \boldsymbol{x} > 0 \ \forall i \neq j \\ \text{(b) Output } h(\boldsymbol{x}) = \texttt{MFAS}(G) \\ \end{array}$

Since $\mathcal{M}(\sigma^*(\boldsymbol{x}))$ is an η -noisy ranking distribution (see Definition 1), we get that the above quantity is at most $\eta < 1/2$. Therefore, each sample (\boldsymbol{x}, y_{ij}) can be viewed as a sample from a distribution \mathcal{D}_{ij} with Gaussian \boldsymbol{x} -marginal, optimal linear classifier sgn $((\boldsymbol{W}_i^* - \boldsymbol{W}_j^*) \cdot \boldsymbol{x})$, and Massart noise η . Hence, we have reduced the task of learning noisy LSFs to a number of $\binom{k}{2}$ sub-problems concerning the learnability of halfspaces in the presence of bounded (Massart) noise.

Step 2: Solving Binary Sub-problems. We can now apply the algorithm MassartLTF for LTFs with Massart noise under standard Gaussian marginals [ZSA20] (for details, see Supp. Material A): for all the pairs of alternatives $1 \le i < j \le k$ with accuracy parameter ϵ' , confidence $\delta' = O(\delta/k^2)$, and a total number of $N = \tilde{\Omega}\left(\frac{d}{\epsilon'(1-2\eta)^6}\log(k/\delta)\right)$ i.i.d. samples from \mathcal{D} , we can obtain a collection of linear classifiers with normal vectors v_{ij} for any i < j. We remark that each one of these halfspaces v_{ij} achieves ϵ disagreement with the ground-truth halfspaces $W_i^* - W_j^*$ with high probability, i.e.,

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\operatorname{sgn}(\boldsymbol{v}_{ij} \cdot \boldsymbol{x}) \neq \operatorname{sgn}((\boldsymbol{W}_i^\star - \boldsymbol{W}_j^\star) \cdot \boldsymbol{x})] \leq \epsilon' \,.$$

Step 3: Ranking Phase. We now have to aggregate the linear classifiers and compute a single 234 sorting function $h : \mathbb{R}^d \to \mathbb{S}_k$. Given an example x, we create the tournament graph G with k nodes 235 that contains a directed edge $e_{i \to j}$ if $v_{ij} \cdot x > 0$. If G is acyclic, we output the induced permutation; 236 otherwise, the graph contains cycles which should be eliminated. In order to output a ranking, we 237 remove cycles from G with an efficient, 3-approximation algorithm for MFAS [ACN08, VZW09]. 238 Hence, the output $h(\boldsymbol{x})$ and the true target $\sigma^{\star}(\boldsymbol{x})$ will have $\mathbf{E}_{\boldsymbol{x}\sim\mathcal{N}_d}[\Delta_{\mathrm{KT}}(h(\boldsymbol{x}),\sigma^{\star}(\boldsymbol{x}))] \leq \epsilon' + 3\epsilon' = 1$ 239 $4\epsilon'$. This last equation indicates why a constant factor approximation algorithm suffices for our 240 purposes – we can always pick $\epsilon' = \epsilon/4$ and complete the proof. For details, see Supp. Material A. 241

242 3.2 Proper Learning Algorithm: Theorem 1

Having obtained the improper learning algorithm, we can now describe our proper Algorithm 2.
Initially, the algorithm starts similarly with the improper learner and obtains a collection of binary
linear classifiers. The crucial idea is the next step: the design of an appropriate convex program which
will efficiently give the matrix *W*. We proceed with the details. For the proof, see Supp. Material A.

Algorithm 2 Proper Learning Algorithm ProperLSF

Input: Training set $T = \{(\boldsymbol{x}^t, \pi^t)\}_{t \in [N]}, \epsilon, \delta \in (0, 1), \eta \in [0, 1/2)$ Output: Linear Sorting function $h : \mathbb{R}^d \to \mathbb{S}_k$, i.e., $h(\cdot) = \sigma_{\boldsymbol{W}}(\cdot)$ for some matrix $\boldsymbol{W} \in \mathbb{R}^{k \times d}$

Compute $(v_{ij})_{1 \le i \le j \le k} = \text{ImproperLSF}(T, \epsilon, \delta, \eta)$	⊳ See Algorithm 1
Setup the CP 1 and compute $W = \text{Ellipsoid}(CP)$	⊳ See Supp. Material A
Ranking Phase: Given $oldsymbol{x} \in \mathbb{R}^d,$ output $h(oldsymbol{x}) = \mathrm{a}$	$\operatorname{rgsort}(\boldsymbol{W}\boldsymbol{x})$

Step 1: Calling Non-proper Learners. As a first step, the algorithm calls Algorithm 1 with parameters ϵ , δ and $\eta \in [0, 1/2)$ and obtains a list of linear classifiers with normal vectors v_{ij} for i < j. Without loss of generality, assume that $||v_{ij}||_2 = 1$.

Step 2: Designing and Solving the CP 1. Our main goal is to find a matrix W whose LSF is close to the true target in KT distance. We show the following lemma that connects the KT distance between two LSFs with the angle metric $d_{\text{angle}}(\cdot, \cdot)$ defined in Eq. (1). The proof can be found in the Supp. Material A.

Lemma 2. For $W, W^* \in \mathbb{R}^{k \times d}$, it holds $\mathbf{E}_{\boldsymbol{x} \sim \mathcal{N}_d}[\Delta_{\mathrm{KT}}(\sigma_{\boldsymbol{W}}(\boldsymbol{x}), \sigma_{\boldsymbol{W}^*}(\boldsymbol{x}))] \leq d_{\mathrm{angle}}(\boldsymbol{W}, \boldsymbol{W}^*)$.

The above lemma states that, for our purposes, it suffices to control the d_{angle} metric between the guess W and the true matrix W^* . It turns out that, given the binary classifiers v_{ij} , we can design a convex program whose solution will satisfy this property. Thinking of the binary classifier v_{ij} as a proxy for $W_i^* - W_j^*$, we want each difference $W_i - W_j$ to have small angle with v_{ij} or equivalently to have large correlation with it, i.e., $(W_i - W_j) \cdot v_{ij} \approx ||W_i - W_j||_2$. To enforce this condition, we can therefore use the second order conic constraint $(W_i - W_j) \cdot v_{ij} \ge (1 - \phi) ||W_i - W_j||_2$. We formulate the following convex program 1 with variable the matrix W:

Find $\boldsymbol{W} \in \mathbb{R}^{k \times d}$, $\|\boldsymbol{W}\|_F \leq 1$, such that $(\boldsymbol{W}_i - \boldsymbol{W}_j) \cdot \boldsymbol{v}_{ij} \geq (1 - \phi) \cdot \|\boldsymbol{W}_i - \boldsymbol{W}_j\|_2$ for any $1 \leq i < j \leq k$, (1)

for some $\phi \in (0, 1)$ to be decided. Intuitively, since any v_{ij} has good correlation with $W_i^* - W_j^*$ (by the guarantees of the improper learning algorithm) and the CP 1 requires that its solution Wsimilarly correlates well with v_{ij} , we expect that $d_{angle}(W, W^*)$ will be small. We show that:

Claim 1. The convex program 1 is feasible and any solution W of 1 satisfies $d_{angle}(W, W^*) \leq \epsilon$.

To see this, note that any solution of CP 1 is a matrix W whose angle metric (see Eq. (1)) with the true matrix is small by an application of the triangle inequality between the angles of $(v_{ij}, W_i - W_j)$ and $(v_{ij}, W_i^* - W_j^*)$ for any $i \neq j$. We next have to deal with the feasibility of CP 1. Our goal is to determine the value of ϕ that makes the CP 1 feasible. For the pair $1 \leq i < j \leq k$, the guess v_{ij} and the true normal vector $W_i^* - W_j^*$ satisfy, with high probability,

$$\Pr_{\boldsymbol{x}\sim\mathcal{D}_{\boldsymbol{x}}}[\operatorname{sgn}(\boldsymbol{v}_{ij}\cdot\boldsymbol{x})\neq\operatorname{sgn}((\boldsymbol{W}_{i}^{\star}-\boldsymbol{W}_{j}^{\star})\cdot\boldsymbol{x})]\leq\epsilon\,.$$
(2)

Under the Gaussian distribution (which is rotationally symmetric), it is well known that the angle $\theta(\boldsymbol{u}, \boldsymbol{v})$ between two vectors $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^d$ is equal to $\pi \cdot \mathbf{Pr}_{\boldsymbol{x} \sim \mathcal{N}_d}[\operatorname{sgn}(\boldsymbol{u} \cdot \boldsymbol{x}) \neq \operatorname{sgn}(\boldsymbol{v} \cdot \boldsymbol{x})]$. Hence, using Eq. (2), we get that the angle between the guess \boldsymbol{v}_{ij} and the true normal vector $\boldsymbol{W}_i^* - \boldsymbol{W}_j^*$ is $\theta(\boldsymbol{W}_i^* - \boldsymbol{W}_j^*, \boldsymbol{v}_{ij}) \leq c\epsilon$. For sufficiently small ϵ , this bound implies that the cosine of the above angle is of order $1 - (c\epsilon)^2$ and so the following inequality will hold (since \boldsymbol{v}_{ij} is unit):

$$(\boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}) \cdot \boldsymbol{v}_{ij} \ge (1 - 2(c\epsilon)^2) \cdot \|\boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}\|_2.$$

Hence, by setting $\phi = 2(c\epsilon)^2$, the convex program 1 with variables $W \in \mathbb{R}^{k \times d}$ will be feasible; since $\|W^*\|_F \leq 1$ comes without loss of generality, W^* will be a solution with probability $1 - \delta$.

Next, we have to control the volume of the feasible region. This is crucial in order to apply the ellipsoid algorithm (for details, see in Supp. Material A) and, hence, solve the convex program. We show the following claim (see Supp. Material A for the proof):

Claim 2. There exists $r \ge 2^{-\text{poly}(d,k,1/\epsilon,\log(1/\delta))}$ so that the feasible set of CP 1 with $\phi = O(\epsilon^2)$ contains a ball (with respect to the Frobenius norm) of radius r.

²⁸³ Critically, the runtime of the ellipsoid algorithm is *logarithmic* in 1/r. So, the ellipsoid runs in time ²⁸⁴ polynomial in the parameters of the problem and outputs the desired matrix W.

4 Learning in top-*r* **Disagreement: Theorem 2**

In this section we show that the proper learning algorithm of Section 3.2 learns noisy LSFs in the top-r286 disagreement metric. We have seen that, with $O(d \log(k)/\epsilon)$ samples, Algorithm 2 of Section 3.2 287 computes a matrix W such that $d_{\text{angle}}(W, W^*) \leq \epsilon$, see Claim 1. Our main contribution is the 288 following lemma that connects the top-r disagreement metric with the geometric distance $d_{angle}(\cdot, \cdot)$, 289 recall Lemma 1. To keep this sketch simple we shall present a sketch of the proof of Lemma 1 for the 290 special case of top-1 classification, which we restate below. The proof of the top-1 case can be found 291 at the Supp. Material B. The detailed proof of the general case (r > 1) can be found in the Supp. 292 Material C. 293

Lemma 3 (Top-1 Disagreement Loss via $d_{angle}(\cdot, \cdot)$). Consider two matrices $U, V \in \mathbb{R}^{k \times d}$ and let \mathcal{N}_d be the standard Gaussian in d dimensions. We have that

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\sigma_1(\boldsymbol{U}\boldsymbol{x}) \neq \sigma_1(\boldsymbol{V}\boldsymbol{x})] \le O\left(k\sqrt{\log k}\right) \ d_{\mathrm{angle}}(\boldsymbol{U}, \boldsymbol{V}) \,.$$

296 We observe that

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[\sigma_1(\boldsymbol{U}\boldsymbol{x})\neq\sigma_1(\boldsymbol{V}\boldsymbol{x})] = \sum_{i\in[k]}\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[\sigma_1(\boldsymbol{U}\boldsymbol{x})=i,\sigma_1(\boldsymbol{V}\boldsymbol{x})\neq i].$$
(1)

We denote by $C_{U}^{(i)} \triangleq \mathbf{1}\{x : \sigma_1(Ux) = i\} = \prod_{j \neq i} \mathbf{1}\{(U_i - U_j) \cdot x \ge 0\}$, i.e., this is the set where the ranking corresponding to U picks i as the top element. Note that $C_{U}^{(i)}$ is the indicator of a homogeneous polyhedral cone since it can be written as the intersection of homogeneous halfspaces. Using these cones we can rewrite the top-1 disagreement of Eq. (1) as

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[\sigma_1(\boldsymbol{U}\boldsymbol{x})\neq\sigma_1(\boldsymbol{V}\boldsymbol{x})] = \sum_{i\in[k]}\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[C_{\boldsymbol{U}}^{(i)}(\boldsymbol{x})=1, C_{\boldsymbol{V}}^{(i)}(\boldsymbol{x})=0].$$
(2)

Hence, our task is to control the mass of the disagreement region of two cones. The next Lemma 4 achieves this task and, combined with Eq. (2) directly gives the conclusion of Lemma 3.

Next we work with two general homogeneous polyhedral cones with set indicators C_1, C_2 :

Lemma 4 (Cone Disagreement). Let $C_1, C_2 : \mathbb{R}^d \mapsto \{0, 1\}$ be homogeneous polyhedral cones defined by the k unit vectors $\mathbf{v}_1, \ldots, \mathbf{v}_k$ and $\mathbf{u}_1, \ldots, \mathbf{u}_k$ respectively. For some universal constant c > 0, it holds that $\mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d}[C_1(\mathbf{x}) \neq C_2(\mathbf{x})] \leq c\sqrt{\log k} \max_{i \in [k]} \theta(\mathbf{v}_i, \mathbf{u}_i)$.

Roadmap of the Proof of Lemma 4: Assume that we rotate one face of the polyhedral cone C_1 by 307 a very small angle θ to obtain the perturbed cone C_2 . At a high-level, we expect the probability of 308 the disagreement region between the new cone C_2 and C_1 to be roughly (this is an underestimation) 309 equal to the size of the perturbation θ times the (Gaussian) surface area of the face of the convex 310 cone that we perturbed. The Gaussian Surface Area (GSA) of a convex set $A \subset \mathbb{R}^d$, is defined 311 as $\Gamma(A) \triangleq \int_{\partial A} \phi_d(\boldsymbol{x}) d\mu(\boldsymbol{x})$, where $d\mu(\boldsymbol{x})$ is the standard surface measure in \mathbb{R}^d and $\phi_d(\boldsymbol{x}) =$ 312 $(2\pi)^{-d/2} \cdot \exp(-\|\boldsymbol{x}\|_2^2/2)$. In fact, in Claim 3 below, we show that the probability of the disagreement 313 between C_1 and C_2 is roughly $O(\theta)\Gamma(F_1)\sqrt{\log(1/\Gamma(F_1)+1)}$, where F_1 is the face of cone C_1 that 314 we rotated. Now, when we perturb all the faces by small angles (all perturbations are at most θ), we 315 can show (via a sequence of triangle inequalities) that the total probability of the disagreement region 316 is bounded above by the perturbation size θ times the sum of the Gaussian surface area of every face 317 (times a logarithmic blow-up factor): 318

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d} [C_1(\boldsymbol{x}) \neq C_2(\boldsymbol{x})] \le O(\theta) \sum_{i=1}^k \Gamma(F_i) \sqrt{\log(1/\Gamma(F_i) + 1)}$$

Surprisingly, for homogeneous convex cones, the above sum cannot grow very fast with k. In fact, we show that it can be at most $O(\sqrt{\log k})$. To prove this, we crucially rely on the following convex geometry result showing that the Gaussian surface area of a homogeneous convex cone is O(1)regardless of the number of its faces k.

Lemma 5 ([Naz03]). Let C be a homogeneous polyhedral cone with k faces F_1, \ldots, F_k . Then C has Gaussian surface area $\Gamma(C) = \sum_{i=1}^k \Gamma(F_i) \le 1$.

Using an inequality similar to the fact that the maximum entropy of a discrete distribution on k elements is at most log k, and, since, from Lemma 5, it holds that $\sum_{i=1}^{k} \Gamma(F_i) \leq 1$, we can show that $\sum_{i=1}^{k} \Gamma(F_i) \sqrt{\log(1/\Gamma(F_i) + 1)} = O(\sqrt{\log k})$. Therefore, with the above lemma we conclude that, if the maximum angle perturbation that we perform on C_1 is θ , then the probability of the disagreement region is $O(\theta)$. We next give the formal proof resulting in the upper bound of $O(\sqrt{\log k} \theta)$ for the disagreement.

Single Face Perturbation Bound: Claim 3: We will use the following notation for the positive orthant indicator $R(z) = \prod_{i=1}^{k} \mathbf{1}\{z_i \ge 0\}$. Notice that the homogeneous polyhedral cone C_1 can be written as $C_1(x) = R(\mathbf{V}x) = R(\mathbf{v}_1 \cdot x, \dots, \mathbf{v}_k \cdot x)$. Claim 3 below shows that the disagreement of two cones that differ on a single normal vector is bounded by above by the Gaussian surface area of a particular face F_1 times a logarithmic blow-up factor $\sqrt{\log(1/\Gamma(F_1) + 1)}$. **Claim 3.** Let $v_1, \ldots, v_k \in \mathbb{R}^d$ and $r \in \mathbb{R}^d$ with $\theta(v_1, r) \leq \theta$ for some sufficiently small $\theta \in (0, \pi/2)$. 137 Let F_1 be the face with $v_1 \cdot x = 0$ of the cone R(Vx) and c > 0 be some universal constant. Then,

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d} \left[R(\boldsymbol{v}_1 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \neq R(\boldsymbol{r} \cdot \boldsymbol{x}, \boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \right] \le c \cdot \theta \cdot \Gamma(F_1) \sqrt{\log \left(\frac{1}{\Gamma(F_1)} + 1\right)} \,.$$

Proof Sketch of Claim 3. Since the constraints $v_2 \cdot x \ge 0, \dots, v_k \cdot x \ge 0$ are common in the two cones, we have that $R(v_1 \cdot x, \dots, v_k \cdot x) \ne R(r \cdot x, v_2 \cdot x, \dots, v_k \cdot x)$ only when the first "halfspaces" disagree, i.e., when $(v_1 \cdot x)(r \cdot x) < 0$. Thus, we have that the LHS probability of Claim 3 is equal to

$$\mathbf{E}_{\boldsymbol{x}\sim\mathcal{N}_d}\left[R(\boldsymbol{v}_2\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})\cdot\mathbf{1}\{(\boldsymbol{v}_1\cdot\boldsymbol{x})(\boldsymbol{r}\cdot\boldsymbol{x})<0\}\right].$$
(3)

This expectation contains two terms: the term $R(v_2 \cdot x, \dots v_k \cdot x)$ that contains the last k - 1common constrains of the two cones and the region where the first two halfspaces disagree, i.e., the set $\{x : (v_1 \cdot x)(r \cdot x) < 0\}$. In order to upper bound this integral in terms of the angle θ , we observe that (for θ sufficiently small) it is not hard to show (see Supp. Material B) that the disagreement region, which is itself a (non-convex) cone, is a subset of the region $\{x : |v_1 \cdot x| \le 2\theta | q \cdot x|\}$, where q the normalized projection of r onto the orthogonal complement of v_1 , i.e., $q = \text{proj}_{v_1^\perp} r / \|\text{proj}_{v_1^\perp} r\|_2$. Therefore, we have that the integral of Eq. (3) is at most

$$\mathbf{E}_{\boldsymbol{x} \sim \mathcal{N}_d} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \ \mathbf{1} \{ |\boldsymbol{v}_1 \cdot \boldsymbol{x}| \leq 2\theta |\boldsymbol{q} \cdot \boldsymbol{x}| \} \right]$$

This is where the definition of the Gaussian surface area appears. In fact, we have to compute the derivative of the above expression (which is a function of θ) with respect to θ and evaluate it at $\theta = 0$. The idea behind this computation is that we can upper bound probability mass of the cone disagreement, i.e., the term $\Pr_{\boldsymbol{x} \sim \mathcal{N}_d} [R(\boldsymbol{v}_1 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \neq R(\boldsymbol{r} \cdot \boldsymbol{x}, \boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x})]$ by its derivative with respect to θ (evaluated at 0) times θ by introducing $o(\theta)$ error. Hence, it suffices to upper bound the value of this derivative at 0, which is:

$$2 \mathop{\mathbf{E}}_{\boldsymbol{x} \sim \mathcal{N}_d} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) | \boldsymbol{q} \cdot \boldsymbol{x} | \, \delta(|\boldsymbol{v}_1 \cdot \boldsymbol{x}|) \right] \,,$$

where δ is the Dirac delta function. Notice that, if we did not have the term $|\mathbf{q} \cdot \mathbf{x}|$, the above expression would be exactly equal to two times the Gaussian surface area of the face with $\mathbf{v}_1 \cdot \mathbf{x} = 0$, i.e., it would be equal to $2\Gamma(F_1)$. We now show that this extra term of $|\mathbf{q} \cdot \mathbf{x}|$ can only increase the above surface integral by at most a logarithmic factor. For some ξ to be decided, we have that

$$\begin{split} \mathbf{E}_{\boldsymbol{x}\sim\mathcal{N}_{d}} \left[R(\boldsymbol{v}_{2}\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_{k}\cdot\boldsymbol{x}) \left| \boldsymbol{q}\cdot\boldsymbol{x} \right| \delta(|\boldsymbol{v}_{1}\cdot\boldsymbol{x}|) \right] &= \int_{\boldsymbol{x}\in F_{1}} \phi_{d}(\boldsymbol{x}) |\boldsymbol{q}\cdot\boldsymbol{x}| d\mu(\boldsymbol{x}) \\ &\leq \int_{\boldsymbol{x}\in F_{1}} \phi_{d}(\boldsymbol{x}) |\boldsymbol{q}\cdot\boldsymbol{x}| \mathbf{1}\{ |\boldsymbol{q}\cdot\boldsymbol{x}| \leq \xi \} d\mu(\boldsymbol{x}) + \int_{\boldsymbol{x}\in F_{1}} \phi_{d}(\boldsymbol{x}) |\boldsymbol{q}\cdot\boldsymbol{x}| \mathbf{1}\{ |\boldsymbol{q}\cdot\boldsymbol{x}| \geq \xi \} d\mu(\boldsymbol{x}) \\ &\leq \xi \int_{\boldsymbol{x}\in F_{1}} \phi_{d}(\boldsymbol{x}) d\mu(\boldsymbol{x}) + \int_{\boldsymbol{x}\in F_{1}} \phi_{d}(\boldsymbol{x}) |\boldsymbol{q}\cdot\boldsymbol{x}| \mathbf{1}\{ |\boldsymbol{q}\cdot\boldsymbol{x}| \geq \xi \} d\mu(\boldsymbol{x}) \,, \end{split}$$

where $d\mu(\boldsymbol{x})$ is the standard surface measure in \mathbb{R}^d . The first integral above is exactly equal to the Gaussian surface area of the face F_1 . To bound from above the second term we can use the next claim showing that not a lot of mass of the face F_1 can concentrate on the region where $|\boldsymbol{q} \cdot \boldsymbol{x}|$ is very large. Its proof relies on standard Gaussian concentration arguments, and is provided in Supp. Material B.

362 **Claim 4.** It holds that
$$\int_{x \in F_1} \phi_d(x) |\mathbf{q} \cdot x| \mathbf{1}\{|\mathbf{q} \cdot x| \ge \xi\} d\mu(x) \le O(\exp(-\xi^2/2))$$
.

³⁶³ Using the above result, we get that

$$\frac{d}{d\theta} \Big(\mathbf{E}_{\boldsymbol{x} \sim \mathcal{N}_d} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \, \mathbf{1} \{ |\boldsymbol{v}_1 \cdot \boldsymbol{x}| \le 2\theta | \boldsymbol{q} \cdot \boldsymbol{x}| \} \right] \Big) \Big|_{\theta=0} \le O(\xi) \, \Gamma(F_1) + O(\exp(-\xi^2/2)) \, .$$

By picking $\xi = \Theta(\sqrt{\log(1 + 1/\Gamma(F_1))})$, the result follows since, up to introducing $o(\theta)$ error, we can bound the term $\operatorname{Pr}_{\boldsymbol{x} \sim \mathcal{N}_d} [R(\boldsymbol{v}_1 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \neq R(\boldsymbol{r} \cdot \boldsymbol{x}, \boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x})]$ by its derivative with respect to θ , evaluated at 0, times θ .

Conclusion. Our work presents the first theoretical guarantees for (linear) LR with noise and settles interesting directions for future work, as mentioned in Section 1. This paper is theoretical and does not have any negative social impact.

370 **References**

371 372 373	[ABHU15]	Pranjal Awasthi, Maria-Florina Balcan, Nika Haghtalab, and Ruth Urner. Efficient learning of linear separators under bounded noise. In <i>Conference on Learning Theory</i> , pages 167–190. PMLR, 2015.
374 375 376	[ABHZ16]	Pranjal Awasthi, Maria-Florina Balcan, Nika Haghtalab, and Hongyang Zhang. Learn- ing and 1-bit compressed sensing under asymmetric noise. In <i>Conference on Learning</i> <i>Theory</i> , pages 152–192. PMLR, 2016.
377 378	[ABSV14]	Pranjal Awasthi, Avrim Blum, Or Sheffet, and Aravindan Vijayaraghavan. Learning mixtures of ranking models. <i>arXiv preprint arXiv:1410.8750</i> , 2014.
379 380	[ACN08]	Nir Ailon, Moses Charikar, and Alantha Newman. Aggregating inconsistent informa- tion: ranking and clustering. <i>Journal of the ACM (JACM)</i> , 55(5):1–27, 2008.
381 382	[AGM17]	Juan A Aledo, José A Gámez, and David Molina. Tackling the supervised label ranking problem by bagging weak learners. <i>Information Fusion</i> , 35:38–50, 2017.
383 384	[AL88]	Dana Angluin and Philip Laird. Learning from noisy examples. <i>Machine Learning</i> , 2(4):343–370, 1988.
385 386	[APA18]	Arpit Agarwal, Prathamesh Patil, and Shivani Agarwal. Accelerated spectral ranking. In <i>International Conference on Machine Learning</i> , pages 70–79. PMLR, 2018.
387 388 389	[BDCBL92]	Shai Ben-David, Nicolò Cesa-Bianchi, and Philip M Long. Characterizations of learnability for classes of $\{O,, n\}$ -valued functions. In <i>Proceedings of the fifth annual workshop on Computational learning theory</i> , pages 333–340, 1992.
390 391 392	[BFFSZ19]	Róbert Busa-Fekete, Dimitris Fotakis, Balázs Szörényi, and Manolis Zampetakis. Optimal learning of mallows block model. In <i>Conference on Learning Theory</i> , pages 529–532. PMLR, 2019.
393	[BH20]	Maria-Florina Balcan and Nika Haghtalab. Noise in classification., 2020.
394 395	[BM09]	Mark Braverman and Elchanan Mossel. Sorting from noisy information. <i>arXiv preprint arXiv:0910.1191</i> , 2009.
396 397	[BT52]	Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. <i>Biometrika</i> , 39(3/4):324–345, 1952.
398 399 400	[BZ17]	Maria-Florina F Balcan and Hongyang Zhang. Sample and computationally efficient learning algorithms under s-concave distributions. <i>Advances in Neural Information Processing Systems</i> , 30, 2017.
401 402	[CDH10]	Weiwei Cheng, Krzysztof Dembczynski, and Eyke Hüllermeier. Label ranking methods based on the plackett-luce model. In <i>ICML</i> , 2010.
403 404	[CH08]	Weiwei Cheng and Eyke Hüllermeier. Instance-based label ranking using the mallows model. In <i>ECCBR Workshops</i> , pages 143–157, 2008.
405 406 407	[CH12]	Weiwei Cheng and Eyke Hüllermeier. Probability estimation for multi-class classifica- tion based on label ranking. In <i>Joint European Conference on Machine Learning and</i> <i>Knowledge Discovery in Databases</i> , pages 83–98. Springer, 2012.
408 409 410	[CHH09]	Weiwei Cheng, Jens Hühn, and Eyke Hüllermeier. Decision tree and instance-based learning for label ranking. In <i>Proceedings of the 26th Annual International Conference on Machine Learning</i> , pages 161–168, 2009.
411 412 413	[CKMY20]	Sitan Chen, Frederic Koehler, Ankur Moitra, and Morris Yau. Classification under misspecification: Halfspaces, generalized linear models, and evolvability. <i>Advances in Neural Information Processing Systems</i> , 33:8391–8403, 2020.
414 415 416	[CKS18]	Stephan Clémençon, Anna Korba, and Eric Sibony. Ranking median regression: Learning to order through local consensus. In <i>Algorithmic Learning Theory</i> , pages 212–245. PMLR, 2018.

417 418 419	[CPS13]	Ioannis Caragiannis, Ariel D Procaccia, and Nisarg Shah. When do noisy votes reveal the truth? In <i>Proceedings of the fourteenth ACM conference on Electronic commerce</i> , pages 143–160, 2013.
420 421 422	[CV20]	Stéphan Clémençon and Robin Vogel. A multiclass classification approach to label ranking. In <i>International Conference on Artificial Intelligence and Statistics</i> , pages 1421–1430. PMLR, 2020. URL: https://arxiv.org/abs/2002.09420.
423 424 425	[Dan16]	Amit Daniely. Complexity theoretic limitations on learning halfspaces. In <i>Proceedings</i> of the forty-eighth annual ACM symposium on Theory of Computing, pages 105–117, 2016.
426 427 428	[DGR+14]	Nemanja Djuric, Mihajlo Grbovic, Vladan Radosavljevic, Narayan Bhamidipati, and Slobodan Vucetic. Non-linear label ranking for large-scale prediction of long-term user interests. In <i>Twenty-Eighth AAAI Conference on Artificial Intelligence</i> , 2014.
429 430 431	[DGT19]	Ilias Diakonikolas, Themis Gouleakis, and Christos Tzamos. Distribution-independent pac learning of halfspaces with massart noise. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
432 433	[DK20]	Ilias Diakonikolas and Daniel M Kane. Hardness of learning halfspaces with massart noise. <i>arXiv preprint arXiv:2012.09720</i> , 2020.
434 435 436	[DKK ⁺ 21]	Ilias Diakonikolas, Daniel M Kane, Vasilis Kontonis, Christos Tzamos, and Nikos Zarifis. Learning general halfspaces with general massart noise under the gaussian distribution. <i>arXiv preprint arXiv:2108.08767</i> , 2021.
437 438 439	[DKM05]	Sanjoy Dasgupta, Adam Tauman Kalai, and Claire Monteleoni. Analysis of perceptron- based active learning. In <i>International conference on computational learning theory</i> , pages 249–263. Springer, 2005.
440 441 442	[DKPZ21]	Ilias Diakonikolas, Daniel M Kane, Thanasis Pittas, and Nikos Zarifis. The optimality of polynomial regression for agnostic learning under gaussian marginals. <i>arXiv</i> preprint arXiv:2102.04401, 2021.
443 444 445	[DKT21]	Ilias Diakonikolas, Daniel Kane, and Christos Tzamos. Forster decomposition and learning halfspaces with noise. <i>Advances in Neural Information Processing Systems</i> , 34, 2021.
446 447 448	[DKTZ20]	Ilias Diakonikolas, Vasilis Kontonis, Christos Tzamos, and Nikos Zarifis. Learn- ing halfspaces with massart noise under structured distributions. In <i>Conference on</i> <i>Learning Theory</i> , pages 1486–1513. PMLR, 2020.
449 450 451	[DKZ20]	Ilias Diakonikolas, Daniel Kane, and Nikos Zarifis. Near-optimal sq lower bounds for agnostically learning halfspaces and relus under gaussian marginals. <i>Advances in Neural Information Processing Systems</i> , 33:13586–13596, 2020.
452 453	[DOS18]	Anindya De, Ryan O'Donnell, and Rocco Servedio. Learning sparse mixtures of rankings from noisy information. <i>arXiv preprint arXiv:1811.01216</i> , 2018.
454 455 456 457	[DSBDSS11]	Amit Daniely, Sivan Sabato, Shai Ben-David, and Shai Shalev-Shwartz. Multiclass learnability and the erm principle. In <i>Proceedings of the 24th Annual Conference on Learning Theory</i> , pages 207–232. JMLR Workshop and Conference Proceedings, 2011.
458 459	[DSM03]	Ofer Dekel, Yoram Singer, and Christopher D Manning. Log-linear models for label ranking. <i>Advances in neural information processing systems</i> , 16:497–504, 2003.
460 461	[DSS14]	Amit Daniely and Shai Shalev-Shwartz. Optimal learners for multiclass problems. In <i>Conference on Learning Theory</i> , pages 287–316. PMLR, 2014.
462 463	[dSSKC17]	Cláudio Rebelo de Sá, Carlos Soares, Arno Knobbe, and Paulo Cortez. Label ranking forests. <i>Expert systems</i> , 34(1):e12166, 2017.

464 465 466 467	[FGKP06]	Vitaly Feldman, Parikshit Gopalan, Subhash Khot, and Ashok Kumar Ponnuswami. New results for learning noisy parities and halfspaces. In 2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06), pages 563–574. IEEE, 2006.
468 469 470	[FHMB08]	Johannes Fürnkranz, Eyke Hüllermeier, Eneldo Loza Mencía, and Klaus Brinker. Multilabel classification via calibrated label ranking. <i>Machine learning</i> , 73(2):133–153, 2008.
471 472	[FKP21]	Dimitris Fotakis, Alkis Kalavasis, and Eleni Psaroudaki. Label ranking through nonparametric regression. <i>arXiv preprint arXiv:2111.02749</i> , 2021.
473 474 475	[FKS21]	Dimitris Fotakis, Alkis Kalavasis, and Konstantinos Stavropoulos. Aggregating in- complete and noisy rankings. In <i>International Conference on Artificial Intelligence</i> <i>and Statistics</i> , pages 2278–2286. PMLR, 2021.
476 477 478	[GDGV13]	Mihajlo Grbovic, Nemanja Djuric, Shengbo Guo, and Slobodan Vucetic. Supervised clustering of label ranking data using label preference information. <i>Machine learning</i> , 93(2-3):191–225, 2013.
479 480 481	[GDV12]	Mihajlo Grbovic, Nemanja Djuric, and Slobodan Vucetic. Learning from pairwise preference data using gaussian mixture model. <i>Preference Learning: Problems and Applications in AI</i> , 33, 2012.
482 483 484	[GGK20]	Surbhi Goel, Aravind Gollakota, and Adam Klivans. Statistical-query lower bounds via functional gradients. <i>Advances in Neural Information Processing Systems</i> , 33:2147–2158, 2020.
485 486	[GR09]	Venkatesan Guruswami and Prasad Raghavendra. Hardness of learning halfspaces with noise. <i>SIAM Journal on Computing</i> , 39(2):742–765, 2009.
487 488 489	[Hau18]	David Haussler. Decision theoretic generalizations of the pac model for neural net and other learning applications. In <i>The Mathematics of Generalization</i> , pages 37–116. CRC Press, 2018.
490 491 492	[HFCB08]	Eyke Hüllermeier, Johannes Fürnkranz, Weiwei Cheng, and Klaus Brinker. Label ranking by learning pairwise preferences. <i>Artificial Intelligence</i> , 172(16):1897–1916, 2008.
493 494 495 496	[HPRZ03]	Sariel Har-Peled, Dan Roth, and Dav Zimak. Constraint classification for multiclass classification and ranking. <i>Advances in neural information processing systems</i> , pages 809–816, 2003. URL: https://proceedings.neurips.cc/paper/2002/file/16026d60ff9b54410b3435b403afd226-Paper.pdf.
497 498 499 500	[HSSVG22]	Daniel Hsu, Clayton Sanford, Rocco Servedio, and Emmanouil-Vasileios Vlatakis-Gkaragkounis. Near-optimal statistical query lower bounds for agnostically learning intersections of halfspaces with gaussian marginals. <i>arXiv preprint arXiv:2202.05096</i> , 2022.
501 502	[Hun04]	David R Hunter. Mm algorithms for generalized bradley-terry models. <i>The annals of statistics</i> , 32(1):384–406, 2004.
503 504	[KCS17]	Anna Korba, Stephan Clémençon, and Eric Sibony. A learning theory of ranking aggregation. In <i>Artificial Intelligence and Statistics</i> , pages 1001–1010. PMLR, 2017.
505 506	[Kea98]	Michael Kearns. Efficient noise-tolerant learning from statistical queries. <i>Journal of the ACM (JACM)</i> , 45(6):983–1006, 1998.
507 508	[KGB18]	Anna Korba, Alexandre Garcia, and Florence d'Alché Buc. A structured prediction approach for label ranking. <i>arXiv preprint arXiv:1807.02374</i> , 2018.
509 510 511	[KMS06]	Claire Kenyon-Mathieu and Warren Schudy. How to rank with few errors–a ptas for weighted feedback arc set on tournaments. In <i>ELECTRONIC COLLOQUIUM ON COMPUTATIONAL COMPLEXITY, REPORT NO. 144 (2006)</i> . Citeseer, 2006.

512 513	[KSS94]	Michael J Kearns, Robert E Schapire, and Linda M Sellie. Toward efficient agnostic learning. <i>Machine Learning</i> , 17(2):115–141, 1994.
514 515	[LB11]	Tyler Lu and Craig Boutilier. Learning mallows models with pairwise preferences. In <i>ICML</i> , 2011.
516 517 518	[LM18]	Allen Liu and Ankur Moitra. Efficiently learning mixtures of mallows models. In 2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS), pages 627–638. IEEE, 2018.
519 520	[LM21]	Allen Liu and Ankur Moitra. Robust voting rules from algorithmic robust statistics. <i>arXiv preprint arXiv:2112.06380</i> , 2021.
521 522	[Luc12]	R Duncan Luce. <i>Individual choice behavior: A theoretical analysis</i> . Courier Corporation, 2012.
523 524	[LV07]	László Lovász and Santosh Vempala. The geometry of logconcave functions and sampling algorithms. <i>Random Structures & Algorithms</i> , 30(3):307–358, 2007.
525	[Mal57]	Colin L Mallows. Non-null ranking models. i. <i>Biometrika</i> , 44(1/2):114–130, 1957.
526 527	[MN06]	Pascal Massart and Élodie Nédélec. Risk bounds for statistical learning. <i>The Annals of Statistics</i> , 34(5):2326–2366, 2006.
528 529 530	[MV19]	Oren Mangoubi and Nisheeth K Vishnoi. Nonconvex sampling with the metropolis- adjusted langevin algorithm. In <i>Conference on Learning Theory</i> , pages 2259–2293. PMLR, 2019.
531 532 533	[MW20]	Cheng Mao and Yihong Wu. Learning mixtures of permutations: Groups of pairwise comparisons and combinatorial method of moments. <i>arXiv preprint arXiv:2009.06784</i> , 2020.
534 535	[Nat89]	Balas K Natarajan. On learning sets and functions. <i>Machine Learning</i> , 4(1):67–97, 1989.
536 537 538	[Naz03]	Fedor Nazarov. On the maximal perimeter of a convex set in \mathbb{R}^n with respect to a gaussian measure. In <i>Geometric aspects of functional analysis</i> , pages 169–187. Springer, 2003.
539 540	[NOS17]	Sahand Negahban, Sewoong Oh, and Devavrat Shah. Rank centrality: Ranking from pairwise comparisons. <i>Operations Research</i> , 65(1):266–287, 2017.
541 542	[NT22]	Rajai Nasser and Stefan Tiegel. Optimal sq lower bounds for learning halfspaces with massart noise. <i>arXiv preprint arXiv:2201.09818</i> , 2022.
543 544	[Pap81]	Christos H Papadimitriou. On the complexity of integer programming. <i>Journal of the ACM (JACM)</i> , 28(4):765–768, 1981.
545 546 547	[RdSRSK15]	Cláudio Rebelo de Sá, Carla Rebelo, Carlos Soares, and Arno Knobbe. Distance-based decision tree algorithms for label ranking. In <i>Portuguese Conference on Artificial Intelligence</i> , pages 525–534. Springer, 2015.
548 549	[RS94]	Ronald L Rivest and Robert Sloan. A formal model of hierarchical concept-learning. <i>Information and Computation</i> , 114(1):88–114, 1994.
550 551	[Sch98]	Alexander Schrijver. <i>Theory of linear and integer programming</i> . John Wiley & Sons, 1998.
552 553	[Slo88]	Robert Sloan. Types of noise in data for concept learning. In <i>Proceedings of the first annual Workshop on Computational Learning Theory</i> , pages 91–96, 1988.
554 555 556	[Slo92]	Robert H Sloan. Corrigendum to types of noise in data for concept learning. In <i>Proceedings of the fifth annual workshop on Computational learning theory</i> , page 450, 1992.

557 558	[Slo96]	Robert H Sloan. Pac learning, noise, and geometry. In <i>Learning and Geometry: Computational Approaches</i> , pages 21–41. Springer, 1996.
559 560	[SS07]	Shai Shalev-Shwartz. <i>Online learning: Theory, algorithms, and applications</i> . PhD thesis, The Hebrew University of Jerusalem, 2007.
561 562	[SSBD14]	Shai Shalev-Shwartz and Shai Ben-David. Understanding machine learning: From theory to algorithms. Cambridge university press, 2014.
563 564	[Vap06]	Vladimir Vapnik. <i>Estimation of dependences based on empirical data</i> . Springer Science & Business Media, 2006.
565 566 567	[VG10]	Shankar Vembu and Thomas Gärtner. Label ranking algorithms: A survey. In <i>Preference learning</i> , pages 45–64. Springer, 2010. URL: https://link.springer.com/chapter/10.1007/978-3-642-14125-6_3.
568 569	[Vis21]	Nisheeth K Vishnoi. <i>Algorithms for convex optimization</i> . Cambridge University Press, 2021.
570 571 572	[VZW09]	Anke Van Zuylen and David P Williamson. Deterministic pivoting algorithms for constrained ranking and clustering problems. <i>Mathematics of Operations Research</i> , 34(3):594–620, 2009.
573 574	[YZ17]	Songbai Yan and Chicheng Zhang. Revisiting perceptron: Efficient and label-optimal learning of halfspaces. <i>Advances in Neural Information Processing Systems</i> , 30, 2017.
575 576 577 578 579	[ZL21]	Chicheng Zhang and Yinan Li. Improved algorithms for efficient active learning halfspaces with massart and tsybakov noise. In Mikhail Belkin and Samory Kpotufe, editors, <i>Proceedings of Thirty Fourth Conference on Learning Theory</i> , volume 134 of <i>Proceedings of Machine Learning Research</i> , pages 4526–4527. PMLR, 15–19 Aug 2021.
580 581 582	[ZLC17]	Yuchen Zhang, Percy Liang, and Moses Charikar. A hitting time analysis of stochastic gradient langevin dynamics. In <i>Conference on Learning Theory</i> , pages 1980–2022. PMLR, 2017.
583 584 585	[ZLGQ14]	Yangming Zhou, Yangguang Liu, Xiao-Zhi Gao, and Guoping Qiu. A label ranking method based on gaussian mixture model. <i>Knowledge-Based Systems</i> , 72:108–113, 2014.
586 587	[ZLY+14]	Yangming Zhou, Yangguang Liu, Jiangang Yang, Xiaoqi He, and Liangliang Liu. A taxonomy of label ranking algorithms. <i>J. Comput.</i> , 9(3):557–565, 2014.
588 589	[ZQ18]	Yangming Zhou and Guoping Qiu. Random forest for label ranking. <i>Expert Systems with Applications</i> , 112:99–109, 2018.
590 591 592 593	[ZSA20]	Chicheng Zhang, Jie Shen, and Pranjal Awasthi. Efficient active learning of sparse halfspaces with arbitrary bounded noise. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 7184–7197. Curran Associates, Inc., 2020.

594 Checklist

1. For all authors... 595 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's 596 contributions and scope? [Yes] 597 (b) Did you describe the limitations of your work? [Yes] We have presented the assump-598 tions of our models and we have explained future research directions in Section 1. 599 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We have 600 discussed that in the Conclusion paragraph at the end of the main body. 601 (d) Have you read the ethics review guidelines and ensured that your paper conforms to 602 them? [Yes] 603

604	2. If you are including theoretical results
605 606	(a) Did you state the full set of assumptions of all theoretical results? [Yes] We have specified the theoretical models we are working on, see e.g., Section 1.1.
607 608	(b) Did you include complete proofs of all theoretical results? [Yes] Due to space con- straints we have included the full proofs of the results in the Supp. Material.
609	3. If you ran experiments
610 611	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [N/A]
612 613	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
614 615	(c) Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [N/A]
616 617	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
618	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
619	(a) If your work uses existing assets, did you cite the creators? [N/A]
620	(b) Did you mention the license of the assets? [N/A]
621 622	(c) Did you include any new assets either in the supplemental material or as a URL? $[N/A]$
623 624	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
625 626	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
627	5. If you used crowdsourcing or conducted research with human subjects
628 629	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
630 631	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
632 633	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

⁶³⁴ A Learning LSFs with Bounded Noise in Kendall's Tau distance

635 A.1 Improperly Learning LSFs with Bounded Noise

We provide an improper learner for LSFs in the presence of bounded noise. We first restate the main result of this section, whose proof relies on a connection between noisy linear label ranking distributions and the Massart noise model.

Theorem 3 (Non-Proper Learning Algorithm). Fix $\eta \in [0, 1/2)$ and $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2. ImproperLSF (Algorithm 1) draws $N = \widetilde{O}\left(\frac{d}{\epsilon(1-2\eta)^6}\log(k/\delta)\right)$ samples from \mathcal{D} , runs in poly $(d, k, 1/\epsilon, \log(1/\delta))$ time and, with probability at least $1 - \delta$, outputs a hypothesis $h : \mathbb{R}^d \to \mathbb{S}_k$ that is ϵ -close in KT distance to the target.

Proof. Assume that the target function is $\sigma^{\star}(\boldsymbol{x}) = \sigma_{\boldsymbol{W}^{\star}}(\boldsymbol{x}) = \operatorname{argsort}(\boldsymbol{W}^{\star}\boldsymbol{x})$ for some unknown matrix $\boldsymbol{W}^{\star} \in \mathbb{R}^{k \times d}$. Consider a collection of N i.i.d. samples from an η -noisy linear label ranking distribution \mathcal{D} (see Definition 2) and let T be the associated training set. For each example $(\boldsymbol{x}, \pi) \in T$, we create a list of $\binom{k}{2}$ binary examples (\boldsymbol{x}, y_{ij}) with $y_{ij} = \operatorname{sgn}(\pi(i) - \pi(j))$ for any $1 \le i < j \le k$, where $\pi(i)$ denotes the position of the element i. Hence, we create the datasets T_{ij} consisting of the binary labeled examples (\boldsymbol{x}, y_{ij}) . We have that

$$\Pr_{(\boldsymbol{x}, \pi) \sim \mathcal{D}} \left[y_{ij} \cdot \operatorname{sgn}((\boldsymbol{W}_i^\star - \boldsymbol{W}_j^\star) \cdot \boldsymbol{x}) < 0 \mid \boldsymbol{x} \right] = \Pr_{\pi \sim \mathcal{M}(\sigma^\star(\boldsymbol{x}))} \left[\pi(i) < \pi(j) \mid \boldsymbol{W}_i^\star \cdot \boldsymbol{x} < \boldsymbol{W}_j^\star \cdot \boldsymbol{x} \right] \,.$$

Since $\mathcal{M}(\sigma^*(\boldsymbol{x}))$ is an η -bounded noise ranking distribution (see Definition 1), we get that

$$\Pr_{\pi \sim \mathcal{M}(\sigma^{\star}(\boldsymbol{x}))} \left[\pi(i) < \pi(j) \mid \sigma^{\star}(\boldsymbol{x})(i) > \sigma^{\star}(\boldsymbol{x})(j) \right] \leq \eta < 1/2$$

where $\sigma^{\star}(\boldsymbol{x})(i)$ denotes the position of the element *i* in the ranking $\sigma^{\star}(\boldsymbol{x})$. Focusing on the training set T_{ij} , we have that the sign y_{ij} is flipped with probability at most η . So, we have reduced the problem to $\binom{k}{2}$ sub-problems concerning the learnability of halfspaces in the presence of Massart noise. The Massart noise model is a special case of Definition 2 where k = 2. Note also that for each training set T_{ij} , the features \boldsymbol{x} have the same distribution. We can now apply the following result for LTFs with Massart noise for the standard Gaussian distribution. Recall that the concept class of homogeneous halfspaces (or linear threshold functions) is $C_{\text{LTF}} = \{h_{\boldsymbol{w}}(\boldsymbol{x}) = \text{sgn}(\boldsymbol{w} \cdot \boldsymbol{x}) : \boldsymbol{w} \in \mathbb{R}^d\}$.

Lemma 6 (Learning Halfspaces with Massart noise [ZSA20]). Fix $\eta \in [0, 1/2)$ and let $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2 with k = 2 (where $C_{\text{LSF}} = C_{\text{LTF}}$). There is a computationally efficient algorithm MassartLTF that draws $m = O(\frac{d \operatorname{polylog}(d)}{\epsilon(1-2\eta)^6} \cdot \log(1/\delta))$ samples from \mathcal{D} , runs in $\operatorname{poly}(m)$ time and outputs a linear threshold function h that is ϵ -close to the target linear threshold function h^* with probability at least $1 - \delta$, i.e., it holds $\operatorname{Pr}_{\boldsymbol{x} \sim \mathcal{N}_d}[h(\boldsymbol{x}) \neq h^*(\boldsymbol{x})] \leq \epsilon$.

We can invoke the algorithm of Lemma 6 for any alternatives $1 \le i < j \le k$ with accuracy $\epsilon' = O(\epsilon)$, $\delta' = O(\delta/k^2)$ and error rate $\eta < 1/2^4$. We remark that Lemma 6 returns a halfspace. Each one of the $\binom{k}{2}$ calls will provide a vector $v_{ij} \in \mathbb{R}^d$ such that, with probability at least $1 - \delta'$, it satisfies

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[\mathrm{sgn}(\boldsymbol{v}_{ij}\cdot\boldsymbol{x})\neq\mathrm{sgn}((\boldsymbol{W}_i^{\star}-\boldsymbol{W}_j^{\star})\cdot\boldsymbol{x})]\leq\epsilon'\,,$$

where the true target halfspace has normal vector $W_i^{\star} - W_j^{\star}$. Moreover, for any i < j, the algorithm requires that the training set T_{ij} is of size

$$|T_{ij}| = \Omega\left(\frac{d}{\epsilon'} \cdot \frac{1}{(1-2\eta)^6} \cdot \log(1/\delta')\right) \,,$$

and, so, a total number of

$$N = \Omega\left(rac{d}{\epsilon} \cdot rac{1}{(1-2\eta)^6} \cdot \log(k/\delta)
ight),$$

1

⁴We can assume that η is known without loss of generality.

samples (\boldsymbol{x}, π) is required from the distribution \mathcal{D} . Given a collection of linear classifiers with normal vectors \boldsymbol{v}_{ij} for any i < j, it remains to aggregate them and compute a sorting function $h : \mathbb{R}^d \to \mathbb{S}_k$. To this end, the estimator h, given an example \boldsymbol{x} , creates the directed complete graph G with k nodes with directed edge $i \to j$ if $\boldsymbol{v}_{ij} \cdot \boldsymbol{x} > 0$. If all the linear classifiers are correct (which occurs with probability $1 - O(\epsilon k^2)$ over \mathcal{D}_x due to the union bound), the graph G is acyclic (since it will match the true directions induced by \boldsymbol{W}^*) and the estimator h outputs the induced permutation. Observe that the KT distance is

$$\frac{1}{\binom{k}{2}} \cdot \mathop{\mathbf{E}}_{\boldsymbol{x} \sim \mathcal{N}_d} \left[\sum_{1 \leq i < j \leq k} \mathbf{1} \{ \operatorname{sgn}(\boldsymbol{v}_{ij} \cdot \boldsymbol{x}) \neq \operatorname{sgn}((\boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}) \cdot \boldsymbol{x}) \} \right] \leq \epsilon' \, .$$

Otherwise, the classifiers are inconsistent and *G* contains cycles. So, the expected number of mistakes in the graph *G* is ϵk^2 . The estimator in order to output a ranking uses a deterministic constant approximation algorithm for the minimum Feedback Arc Set [ACN08] in order to remove the cycles. For an overview of this fundamental line of research, we refer to [ACN08, VZW09, KMS06].

Lemma 7 (3-Approximation Algorithm for mimimum FAS (see [VZW09, ACN08])). *There is a deterministic algorithm MFAS for the minimum Feedback Arc Set on unweighted tournaments with k vertices that outputs orderings with cost less than* $3 \cdot \text{OPT}$. *The running time is* poly(k).

In the above, OPT is the minimum number of flips the algorithm should perform. With input the cyclic directed graph G induced by the estimated linear classifiers, the algorithm of Lemma 7 computes, in poly(k) time, a 3-approximation of the optimal solution (i.e., instead of correcting ϵ_0 directed edges, the algorithm will provide a directed acyclic graph with $3\epsilon_0$ changed edges). Hence, for the hypothesis $h : \mathbb{R}^d \to \mathbb{S}_k$, where h(x) is the output of the minimum FAS approximation algorithm with input G (G depends on the input x, the randomness of the samples and the internal randomness of the $\binom{k}{2}$ calls of the Massart linear classifiers), and the target function $\sigma^*(x)$, we have that

$$\mathop{\mathbf{E}}_{\boldsymbol{x}\sim\mathcal{N}_d}[\Delta_{KT}(h(\boldsymbol{x}),\sigma^{\star}(\boldsymbol{x}))] \leq (\epsilon'+3\epsilon') = 4\epsilon'$$

which completes the proof, by setting $\epsilon' = \epsilon/4$.

Remark 1. Consider the following variant of the above procedure: compute the $O(k^2)$ linear classifiers with accuracy $\epsilon' = \epsilon/k^2$: If the induced directed graph is acyclic, output the ranking; otherwise, output a random permutation. With probability ϵ , the KT distance will be of order k^2 . Hence, one has to draw in total $O(k^4d/\epsilon)$ samples to make the expected KT distance roughly $O(\epsilon)$. The algorithm of Theorem 3 improves on this approach.

698 A.2 The Proof of Theorem 1: Properly Learning LSFs with Bounded Noise

699 We first restate the main result of this section.

Theorem 4 (Proper Learning Algorithm). Fix $\eta \in [0, 1/2)$ and $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2. ProperLSF (Algorithm 2) draws $N = \widetilde{O}\left(\frac{d}{\epsilon(1-2\eta)^6}\log(k/\delta)\right)$ samples from \mathcal{D} , runs in poly $(d, k, 1/\epsilon, \log(1/\delta))$ time and, with probability at least $1 - \delta$, outputs a Linear Sorting function $h : \mathbb{R}^d \to \mathbb{S}_k$ that is ϵ -close in KT distance to the target.

We are now ready to provide the proof of our efficient proper learning algorithm for the class of Linear Sorting functions in the presence of bounded noise with respect to the standard Gaussian probability measure.

Proof. As a first step, the algorithm calls the improper learning algorithm ImproperLSF (Algorithm 1) with parameters ϵ , δ and $\eta < 1/2$ and obtains a list of linear classifiers with normal vectors v_{ij} for i < j. The utility of this step implies that, with probability at least $1 - \delta$, each one of the classifiers ϵ -learns the associated true halfspace, i.e., it holds

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\operatorname{sgn}(\boldsymbol{v}_{ij} \cdot \boldsymbol{x}) \neq \operatorname{sgn}((\boldsymbol{W}_i^\star - \boldsymbol{W}_j^\star) \cdot \boldsymbol{x})] \leq \epsilon \,,$$

- where W^* is the matrix of the target Linear Sorting function. Without loss of generality, assume that
- 713 $\|v_{ij}\|_2 = 1$. In order to make the learner proper, it suffices to solve the following convex program on
- 714 W:

SU

Find
$$\boldsymbol{W} \in \mathbb{R}^{k \times d}$$
, (1)

ch that
$$(\boldsymbol{W}_i - \boldsymbol{W}_j) \cdot \boldsymbol{v}_{ij} \ge (1 - \phi) \cdot \|\boldsymbol{W}_i - \boldsymbol{W}_j\|_2$$
 for any $1 \le i < j \le k$, (CP) (2)
 $\|\boldsymbol{W}\|_F \le 1$, (3)

for some $\phi \in (0, 1)$ to be decided. The main key ideas are summarized in the next claim.

- **Claim 5.** The following properties hold true for $\phi = O(\epsilon^2)$ with probability at least 1δ .
- *1. The convex program 1 is feasible.*
- 718 2. Any solution of the convex program 1 induces an LSF that is ϵ -close in KT distance to the 719 true target $\sigma_{W^*}(\cdot)$.
- 720 3. The feasible set of the convex program 1 contains a ball of radius $r = 2^{-\text{poly}(d,k,1/\epsilon,\log(1/\delta))}$ 721 and is contained in a ball of radius 1. Both balls are with respect to the Frobenius norm.
- *4. The convex program 1 can be solved in time* $poly(d, k, 1/\epsilon, log(1/\delta))$ *using the ellipsoid algorithm.*

Proof of Item 1. First, we can choose the error ϕ so that this convex program is feasible. Let us set $W = W^*$, where W^* is the underlying matrix of the target Linear Sorting function σ^* with $\sigma^*(x) = \operatorname{argsort}(W^*x)$. Recall that, by the guarantees of the improper learning algorithm, for the pair $1 \le i < j \le k$, it holds

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[\operatorname{sgn}(\boldsymbol{v}_{ij}\cdot\boldsymbol{x})\neq\operatorname{sgn}((\boldsymbol{W}_i^{\star}-\boldsymbol{W}_j^{\star})\cdot\boldsymbol{x})]\leq\epsilon\,.$$
(4)

Since the standard Gaussian is rotationally symmetric, the angle $\theta(u, v)$ between two vectors $u, v \in \mathbb{R}^d$ is equal to $\pi \cdot \Pr_{x \sim \mathcal{N}_d}[\operatorname{sgn}(u \cdot x) \neq \operatorname{sgn}(v \cdot x)]$. Hence, using this observation and Equation (4), we get that the angle between the guess vector v_{ij} and the true normal vector $W_i^{\star} - W_j^{\star}$ is

$$heta(oldsymbol{W}_i^\star-oldsymbol{W}_j^\star,oldsymbol{v}_{ij})\leq c\cdot\epsilon\,,$$

for some constant c > 0. For sufficiently small ϵ , this bound implies that the cosine of the above angle is of order $1 - (c\epsilon)^2$ and so the following inequality will hold

$$(\boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}) \cdot \boldsymbol{v}_{ij} \geq (1 - 2(c\epsilon)^2) \cdot \|\boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}\|_2,$$

since v_{ij} is unit. Hence, by setting $\phi = 2(c\epsilon)^2$, the convex program with variables $W \in \mathbb{R}^{k \times d}$ will be feasible; W^* will be a solution with probability $1 - \delta$, where the randomness is over the output of the algorithm dealing with the Massart linear classifiers. Note that we can assume that $||W^*||_F \le 1$ without loss of generality, since we can divide each row with the Frobenius norm.

Proof of Item 2. Let \widetilde{W} be a solution of the convex program. We will make use of the observation that the angle between two vectors is equal to the disagreement of the associated linear threshold functions with respect to the standard normal times π . Observe that any solution \widetilde{W} to the convex program will satisfy that

$$(\forall i, j) \quad \theta(\boldsymbol{v}_{ij}, \boldsymbol{W}_i - \boldsymbol{W}_j) \leq O(\sqrt{\phi}) = c\epsilon.$$

742 and

$$(\forall i, j) \quad \theta(\boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}, \boldsymbol{v}_{ij}) \leq \epsilon.$$

743 This implies that

$$d_{\text{angle}}(\boldsymbol{W}^{\star}, \boldsymbol{W}) \leq c' \epsilon$$

Claim 6. For the matrices $W, W^* \in \mathbb{R}^{k \times d}$, it holds that

$$\mathbf{E}_{\boldsymbol{x} \sim \mathcal{N}_d} [\Delta_{\mathrm{KT}}(\sigma_{\boldsymbol{W}}(\boldsymbol{x}), \sigma_{\boldsymbol{W}^{\star}}(\boldsymbol{x}))] \leq d_{\mathrm{angle}}(\boldsymbol{W}, \boldsymbol{W}^{\star}).$$

745 *Proof.* We have that

$$\begin{split} \mathbf{E}_{\boldsymbol{x}\sim\mathcal{N}_{d}}[\Delta_{\mathrm{KT}}(\sigma_{\boldsymbol{W}}(\boldsymbol{x}),\sigma_{\boldsymbol{W}^{\star}}(\boldsymbol{x}))] &= \frac{1}{\binom{k}{2}} \cdot \mathbf{E}_{\boldsymbol{x}\sim\mathcal{N}_{d}}[\sum_{1 \leq i < j \leq k} \mathbf{1}\{((\boldsymbol{W}_{i}-\boldsymbol{W}_{j}) \cdot \boldsymbol{x}) \ ((\boldsymbol{W}_{i}^{\star}-\boldsymbol{W}_{j}^{\star}) \cdot \boldsymbol{x}) < 0\} \\ &= \frac{1}{\binom{k}{2}} \cdot \sum_{1 \leq i < j \leq k} \mathbf{Pr}_{\boldsymbol{x}\sim\mathcal{N}_{d}}[\mathrm{sgn}(\boldsymbol{W}_{i}-\boldsymbol{W}_{j}) \cdot \boldsymbol{x}) \neq \mathrm{sgn}((\boldsymbol{W}_{i}^{\star}-\boldsymbol{W}_{j}^{\star}) \cdot \boldsymbol{x})] \\ &= \frac{1}{\pi} \max_{i,j} \theta(\boldsymbol{W}_{i}-\boldsymbol{W}_{j}, \boldsymbol{W}_{i}^{\star}-\boldsymbol{W}_{j}^{\star}) \\ &\leq d_{\mathrm{angle}}(\boldsymbol{W}, \boldsymbol{W}^{\star}). \end{split}$$

746

Using the above claim, we get an expected KT distance bound of order $O(\epsilon)$. This gives the desired result.

- 749 **Proof of Item 3.** We will make use of the next lemma.
- **Lemma 8.** Fix $\epsilon, \delta \in (0, 1)$. Let $W^* \in \mathbb{R}^{k \times d}$ be the true parameter matrix. There exists a matrix $\widetilde{W}^* \in \mathbb{R}^{k \times d}$ such that, with probability at least $1 - \delta$:

•
$$\mathbf{Pr}_{\boldsymbol{x} \sim \mathcal{N}_d}[\operatorname{sgn}((\boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}) \cdot \boldsymbol{x}) \neq \operatorname{sgn}((\widetilde{\boldsymbol{W}}_i^{\star} - \widetilde{\boldsymbol{W}}_j^{\star}) \cdot \boldsymbol{x})] \leq \epsilon \text{ for all } i \neq j, and,$$

$$\bullet \|\widetilde{W}_i^{\star} - \widetilde{W}_j^{\star}\|_2 \geq 2^{-\operatorname{poly}(d,k,1/\epsilon,\log(1/\delta))} \text{ for any } i \neq j$$

Proof of Lemma 8. The above lemma is a result of the next Appendix A.2.1. In particular, it is a direct implication of Lemma 10 and Corollary 1.

756 Note that the above lemma implies that

$$(\forall i,j) \quad \Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\operatorname{sgn}(\boldsymbol{v}_{ij} \cdot \boldsymbol{x}) \neq \operatorname{sgn}((\widetilde{\boldsymbol{W}}_i^\star - \widetilde{\boldsymbol{W}}_j^\star) \cdot \boldsymbol{x})] \leq 2\epsilon \,,$$

with probability at least $1 - 2\delta$. Hence, up to constants, the analysis concerning the feasibility of the true matrix W^* (see Item 1) will still hold for \widetilde{W}^* . From now on we can work with this matrix \widetilde{W}^* which enjoys the "well-conditionedness" property of the second item of the lemma.

We will use the above lemma in order to prove Item 3 which controls the volume of the feasible region: it states that there exist 0 < r < R so that the feasible region of the convex program contains a ball of radius r and is contained in a ball of radius R (where the balls are with respect to the Frobenius norm). Moreover, $r = 2^{-\text{poly}(d,k,1/\epsilon,\log(1/\delta))}$ and R = 1.

For the chosen $\phi \in (0, 1)$, the feasible set contains matrices $\mathbf{W} \in \mathbb{R}^{k \times d}$ that satisfy $\|\mathbf{W} - \widetilde{\mathbf{W}}^{\star}\|_F \leq 2r, r$ to be decided. For any $i \neq j$, we have that the following properties hold:

3.
$$\|\boldsymbol{W} - \widetilde{\boldsymbol{W}}^{\star}\|_{F} \leq 2r$$
 which implies that $\|\boldsymbol{W}_{i} - \widetilde{\boldsymbol{W}}_{i}^{\star}\|_{2} \leq 2r$ for any $i \in [k]$ (ball around feasible point).

770 4. $\|\boldsymbol{v}_{ij}\|_2 = 1.$

Our goal is to prove that for a matrix in the above ball it holds $(W_i - W_j) \cdot v_{ij} \ge (1 - \phi) \|W_i - W_j\|_2$. We have that

$$\begin{split} (\widetilde{\boldsymbol{W}}_{i}^{\star} - \widetilde{\boldsymbol{W}}_{j}^{\star}) \cdot \boldsymbol{v}_{ij} &= (\widetilde{\boldsymbol{W}}_{i}^{\star} - \boldsymbol{W}_{i}) \cdot \boldsymbol{v}_{ij} + (\boldsymbol{W}_{j} - \widetilde{\boldsymbol{W}}_{j}^{\star}) \cdot \boldsymbol{v}_{ij} + (\boldsymbol{W}_{i} - \boldsymbol{W}_{j}) \cdot \boldsymbol{v}_{ij} \\ &\leq \|\widetilde{\boldsymbol{W}}_{i}^{\star} - \boldsymbol{W}_{i}\|_{2} + \|\boldsymbol{W}_{j} - \widetilde{\boldsymbol{W}}_{j}^{\star}\|_{2} + (\boldsymbol{W}_{i} - \boldsymbol{W}_{j}) \cdot \boldsymbol{v}_{ij} \\ &\leq 4r + (\boldsymbol{W}_{i} - \boldsymbol{W}_{j}) \cdot \boldsymbol{v}_{ij} \,. \end{split}$$

773 More to that

$$\begin{split} \|\boldsymbol{W}_{i} - \boldsymbol{W}_{j}\|_{2} &= \|\boldsymbol{W}_{i} - \widetilde{\boldsymbol{W}}_{i}^{\star} + \widetilde{\boldsymbol{W}}_{i}^{\star} - \widetilde{\boldsymbol{W}}_{j}^{\star} + \widetilde{\boldsymbol{W}}_{j}^{\star} - \boldsymbol{W}_{j}\|_{2} \\ &\leq \|\boldsymbol{W}_{i} - \widetilde{\boldsymbol{W}}_{i}^{\star}\|_{2} + \|\widetilde{\boldsymbol{W}}_{i}^{\star} - \widetilde{\boldsymbol{W}}_{j}^{\star}\|_{2} + \|\widetilde{\boldsymbol{W}}_{j}^{\star} - \boldsymbol{W}_{j}\|_{2} \\ &\leq 4r + \|\widetilde{\boldsymbol{W}}_{i}^{\star} - \widetilde{\boldsymbol{W}}_{j}^{\star}\|_{2}, \end{split}$$

and similarly: $\|\mathbf{W}_i - \mathbf{W}_j\|_2 \ge \|\widetilde{\mathbf{W}}_i^{\star} - \widetilde{\mathbf{W}}_j^{\star}\|_2 - 4r.$

775 Combining the above inequalities, we get that

$$(\mathbf{W}_{i} - \mathbf{W}_{j}) \cdot \mathbf{v}_{ij} \ge (\widetilde{\mathbf{W}}_{i}^{\star} - \widetilde{\mathbf{W}}_{j}^{\star}) \cdot \mathbf{v}_{ij} - 4r$$

$$\ge (1 - \phi) \|\widetilde{\mathbf{W}}_{i}^{\star} - \widetilde{\mathbf{W}}_{j}^{\star}\|_{2} - 4r$$

$$\ge (1 - \phi) (\|\mathbf{W}_{i} - \mathbf{W}_{j}\|_{2} - 4r) - 4r$$

$$= (1 - \phi) \|\mathbf{W}_{i} - \mathbf{W}_{j}\|_{2} - 8r.$$

We pick r sufficiently small and of order $2^{-\text{poly}(d,k,1/\epsilon,\log(1/\delta))}$ and get that W is a feasible solution of the convex program. Moreover, we can select R = 1 since $\|\widetilde{W}^{\star}\|_{F} = 1$ without loss of generality, since we can normalize the row differences of \widetilde{W}^{\star} with the norm $\|\widetilde{W}^{\star}\|_{F}$.

Proof of Item 4. We apply the ellipsoid algorithm in order to solve the convex program 1 and compute a matrix $\widetilde{W} \in \mathbb{R}^{k \times d}$. The algorithm ProperLSF outputs the linear sorting function $h(\cdot) = \sigma_{\widetilde{W}}(\cdot)$.

Lemma 9 (Efficiency of the Ellipsoid Algorithm [Vis21]). Suppose that $P \subseteq \mathbb{R}^d$ is a full-dimensional polytope that is contained in a d-dimensional Euclidean ball of radius R > 0 and contains a ddimensional Euclidean ball of radius r > 0. Then, the ellipsoid method outputs a point $\tilde{x} \in P$ after $O(d^2 \log(R/r))$ iterations. Moreover, every iteration can be implemented in $O(d^2 + T_{sep})$ time, where T_{sep} is the time required to answer a single query by the separation oracle.

Assume that Item 3 holds true. Then the algorithm can be used with $r = 2^{-\text{poly}(d,k,1/\epsilon,\log(1/\delta))}$ and R = 1. Hence, the ellipsoid algorithm will provide in time $\text{poly}(d,k,1/\epsilon,\log(1/\delta))$ a point \widetilde{W} that lies in the feasible region of the convex program 1⁵.

790

\$

Remark 2. We remark that both the improper (Algorithm 1) and the proper (Algorithm 2) learning algorithms hold for the more general case where the *x*-marginal lies in the class of isotropic log-concave distributions [LV07]: A distribution \mathcal{D}_x lies inside the class of isotropic log-concave distributions \mathcal{F}_{LC} over \mathbb{R}^d if \mathcal{D}_x has a probability density function f over \mathbb{R}^d such that $\log f$ is concave, its mean is zero, and its covariance is identity, i.e., $\mathbf{E}_{\boldsymbol{x}\sim\mathcal{D}_x}[\boldsymbol{x}\boldsymbol{x}^{\top}] = \boldsymbol{I}$.

796 A.2.1 The proof of Lemma 8

797 We provide the following result.

Lemma 10. Fix $\epsilon, \delta \in (0, 1)$. Let $W^* \in \mathbb{R}^{k \times d}$ be the true parameter matrix. There exists a matrix $W \in \mathbb{R}^{k \times d}$ such that, with probability at least $1 - \delta$:

•
$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\operatorname{sgn}((\boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}) \cdot \boldsymbol{x}) \neq \operatorname{sgn}((\boldsymbol{W}_i - \boldsymbol{W}_j) \cdot \boldsymbol{x})] \leq \epsilon \text{ for all } i \neq j, and,$$

• The bit complexity of
$$W$$
 is $poly(k, d, 1/\epsilon, log(1/\delta))$

Proof. The matrix W will be the output of a linear program that can be used to learn the LSF $\sigma_{W^*}(\cdot)$ in the noiseless setting.

⁵We remark that the runtime will also depend on the time required to answer a single query by the separation oracle. We assume that this time is polynomial in the parameters of our problem and we opt not to track these details in this work.

- Consider the unit sphere S^{d-1} and a δ_0 -cover of the unit sphere with parameter $\delta_0 > 0$ to be decided. For any sample $(\boldsymbol{x}, \pi) \sim \mathcal{D}$ of the 0-noisy linear label ranking distribution, i.e., $\boldsymbol{x} \sim \mathcal{N}_d$ and $\pi = \sigma_{\boldsymbol{W}^{\star}}(\boldsymbol{x})$, we consider the rounded sample $(\tilde{\boldsymbol{x}}, \pi)$ where $\tilde{\boldsymbol{x}}$ is obtained by first projecting $\boldsymbol{x} \in \mathbb{R}^d$
- to S^{d-1} and then by obtaining the closest point of \hat{x} in the cover. The cover's size is $O(1/\delta_0)^d$.
- Let us fix $1 \le i < j \le k$ and set $y_{ij} = \text{sgn}(\pi(i) \pi(j))$. For a training set $\{(\boldsymbol{x}^{(t)}, \pi^{(t)})\}_{t \in [N]}$ of size N, we create the following linear system L_{ij} with variables $\boldsymbol{W} \in \mathbb{R}^{k \times d}$:

$$y_{ij}^{(t)} \left(\boldsymbol{W}_i - \boldsymbol{W}_j \right) \cdot \widetilde{\boldsymbol{x}}^{(t)} \ge 0, \ t \in [N] \quad (\mathbf{L}_{ij})$$

- Consider the concatenation of the linear systems $L = \bigcup_{i < j} L_{ij}$. The number of equations in the linear system of equations L is $N \cdot {k \choose 2}$.
- We first have to show that, with high probability, the system L is feasible, i.e., there exists W that satisfies the system's equations. Note that if we replace $\tilde{x}^{(t)}$ with the original points $x^{(t)}$, the true matrix W^* is a solution to the system. We now have to study the rounded linear system.
- 815 **Claim 7.** The (rounded) linear system L is feasible with high probability.
- *Proof.* In order to show the feasibility of L, we will use the anti-concentration properties of the Gaussian.

Fact 1 ([DKM05]). Let \mathcal{P} be the standard normal distribution over \mathbb{R}^d . For any fixed unit vector $a \in \mathbb{R}^d$ and any $\gamma \leq 1$,

$$\gamma/4 \leq \Pr_{\boldsymbol{x} \sim \mathcal{P}} \left[|\boldsymbol{a} \cdot \frac{\boldsymbol{x}}{\|\boldsymbol{x}\|_2}| \leq \frac{\gamma}{\sqrt{d}} \right] \leq \gamma.$$

Let us focus on the pair $1 \le i < j \le k$. We first observe that scaling all samples to lie on the unit sphere does not affect the feasibility of the system. It suffices to focus on that single halfspace with normal vector $v_{ij} = W_i^* - W_j^* \in \mathbb{R}^d$ and consider the probability of the event that the collection of the N rounded points $\{\tilde{x}^{(t)}\}_t$ with labels $\{y_{ij}^{(t)}\}_t$, that come from N Gaussian vectors $\{x^{(t)}\}_t$ which are linearly separable (with labels $\{y_{ij}^{(t)}\}_t$), becomes non-linearly separable. For this it suffices to control the probability that the rounding procedure flips the label of the data point. Using the union bound, we have that, if the rounding has accuracy δ_0 , the described bad event has probability

$$\Pr_{\boldsymbol{x}^{(1)},\dots,\boldsymbol{x}^{(N)}\sim\mathcal{N}_d}[\exists t\in[N]:\operatorname{sgn}(\boldsymbol{v}_{ij}\cdot\widetilde{\boldsymbol{x}}^{(t)})\neq\operatorname{sgn}(\boldsymbol{v}_{ij}\cdot\boldsymbol{x}^{(t)})]\leq N\cdot\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[|\boldsymbol{v}_{ij}\cdot\boldsymbol{x}/\|\boldsymbol{x}\|_2|\leq 2\delta_0]\leq N\cdot O(\delta_0\sqrt{d})$$

where we remark that the first event is scale invariant and so we can assume that the normal vector is unit, the first inequality follows from the fact that it suffices to control the mass assigned to a strip of width $2\delta_0$ (due to the discretization) and the second inequality follows from Fact 1. We now have to select the discretization. Let $\delta \in (0, 1)$. By choosing $\delta_0 = O(\frac{\delta}{N\sqrt{dk^2}})$, the bad event for all the pairs i < j occurs with probability at most δ , i.e., with probability at least $1 - \delta$, each one of the N drawn i.i.d. samples does not fall in any one of the $\binom{k}{2}$ "bad" strips.

We can now consider the case that the system L is feasible (with the target matrix W^* being a feasible point) that occurs with probability $1 - \delta$. The class of homogenous halfspaces in ddimensions has VC dimension d; therefore, the sample complexity of learning halfspaces using ERM is $O((d + \log(1/\delta))/\epsilon)$. Moreover, in the realizable case, we can implement the ERM using e.g., linear programming and find a solution in $poly(d, 1/\epsilon, \log(1/\delta))$ time. We next focus on the quality of the solution which will give the desired sample complexity.

Claim 8. Assume that the algorithm draws $N = \widetilde{O}(\frac{d + \log(k/\delta)}{\epsilon})$ i.i.d. samples of the form (\boldsymbol{x}, π) with $\boldsymbol{x} \sim \mathcal{N}_d$ and $\pi = \sigma_{\boldsymbol{W}^{\star}}(\boldsymbol{x})$. For any $i \neq j$ and with probability at least $1 - 2\delta$, the solution \boldsymbol{W} of the linear system L satisfies

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\operatorname{sgn}((\boldsymbol{W}_i^{\star} - \boldsymbol{W}_j^{\star}) \cdot \boldsymbol{x}) \neq \operatorname{sgn}((\boldsymbol{W}_i - \boldsymbol{W}_j) \cdot \boldsymbol{x})] \leq \epsilon \,.$$

Proof. Since the matrix W satisfies the sub-system L_{ij} , the result follows using a union bound on the events that (i) the linear system is feasible and (ii) the ERM is a successful PAC learner. **Claim 9.** Consider the solution W of the linear system. Then, W has bounded bit complexity of order poly $(d, k, 1/\epsilon, \log(1/\delta))$.

Proof. We will make use of the following result that relates the size of the input and the output of a linear program using Cramer's rule.

Lemma 11 ([Sch98, Pap81]). Let $A \in \mathbb{Z}^{m \times n}$, $b \in \mathbb{Z}^{m}$, $c \in \mathbb{Z}^{n}$. Consider a linear program min $c \cdot x$ subject to $Ax \leq b$ and $x \geq 0$. Let U be the maximum size of A_{ij} , b_i , c_j . The output of the linear program has size $O(m(nU + n \log(n)))$ bits.

We will apply the above lemma (which holds even by dropping the constraint $x \ge 0$) to our setting 851 where $Aw \ge 0$ where $w = (W_i)_{i \in [k]} \in \mathbb{Q}^{kd}$, i.e., w is the vectorization of the matrix W. Moreover, 852 A is the matrix containing the N (rounded) Gaussian samples $\widetilde{x}^{(t)}$. We have that the matrix A has 853 dimension $N\binom{k}{2} \times kd$ and each entry A_{ij} is an integer and has size at most U = poly(d, k) (since 854 the samples are rounded on the δ_0 -cover of the sphere. Recall that the labels $y_{ij}^{(t)} \in \{-1, +1\}$ and 855 $\widetilde{x}^{(t)}$ lie in the unit sphere. In particular, each row of the matrix A has 2d non-zero entries and is 856 associated with a tuple (i, j, t) for $1 \le i < j \le k$ and $t \in [N]$. Then, it holds that the output has 857 size at most $O(Nk^2(dU + dk \log(dk)))$ bits. So, we get that the output W can be described using 858 at most $poly(d, k, 1/\epsilon, U, log(1/\delta)) = poly(d, k, 1/\epsilon, log(1/\delta))$ bits (due to the size of the entries of 859 the matrix A). 860 П

⁸⁶¹ Combining the above claims, we conclude the proof.

As a corollary of the bounded bit complexity, we obtain the following key result.

Corollary 1. Let $\epsilon > 0$. Assume that $\mathbf{W} \in \mathbb{R}^{k \times d}$ has bit complexity at most $\operatorname{poly}(d, k, 1/\epsilon, \log(1/\delta))$.

864 Then, for any $i, j \in [k]$ with $i \neq j$, it holds that $\|\mathbf{W}_i - \mathbf{W}_j\|_2 > 2^{-\operatorname{poly}(d,k,1/\epsilon,\log(1/\delta))}$.

Proof. First, we can assume that $W_i \neq W_j$ for any $i \neq j$; in case of equal rows, we obtain a low-dimensional instance. Then, since any vector W_i has bounded bit complexity, we have that the difference of any two such vectors, provided that it is non-zero, has a lower bound in its norm, i.e., $\|W_i - W_j\|_2 > 2^{-\text{poly}(d,k,1/\epsilon,\log(1/\delta))}$ for any $i, j \in [k]$.

B69 B Learning in Top-1 Disagreement from Label Rankings

Let us set $\sigma_1(Wx) = \operatorname{argmax}_{i \in [k]} W_i \cdot x$ for $x \in \mathbb{R}^d$. The main result of this section follows.

Theorem 5 (Proper Top-1 Learning Algorithm). Fix $\eta \in [0, 1/2)$ and $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2. There exists an algorithm that draws $N = O\left(\frac{dk\sqrt{\log k}}{\epsilon(1-2\eta)^6}\log(k/\delta)\right)$ samples from \mathcal{D} , runs in poly(N) time and, with probability at least $1 - \delta$, outputs a Linear Sorting function $h : \mathbb{R}^d \to \mathbb{S}_k$ that is ϵ -close in top-1

probability at least $1 - \delta$, outputs a Linear Sorting function $h : \mathbb{R}^a \to \mathbb{S}_k$ that is ϵ -close in top-1 disagreement to the target.

876 *Proof.* Note that the MassartLTF algorithm (see Lemma 6) has the guarantee that it returns a vector 877 w so that

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\operatorname{sgn}(\boldsymbol{w} \cdot \boldsymbol{x}) \neq \operatorname{sgn}(\boldsymbol{w}^{\star} \cdot \boldsymbol{x})] \leq \epsilon \,,$$

with probability $1 - \delta$, where w^* is the target normal vector. Since the above misclassification 878 probability with respect to \mathcal{N}_d is directly connected with the angle $\theta(w, w^*)$, we get that we can 879 control the angle between w and w^* efficiently. Moreover, in our setting, for a matrix $W \in \mathbb{R}^{k \times d}$, 880 there exist $\binom{k}{2}$ homogeneous halfspaces with normal vectors $W_i - W_j$ and so we can control the 881 angles $\theta(W_i - W_j, W_i^{\star} - W_j^{\star})$. In order to deduce the sample complexity bound of Theorem 5, 882 we show the next lemma which essentially bounds the top-1 misclassification error using the angles 883 of these $O(k^2)$ halfspaces. We apply Lemma 12 with U = W and $V = W^*$ and so we can take 884 $\epsilon' = \epsilon/(k\sqrt{\log k})$ and invoke the proper learning algorithm of Algorithm 2. This completes the 885 proof. \square 886

- ⁸⁸⁷ We continue with the proof of our key lemma.
- **Lemma 12** (Misclassification Error). Consider two matrices $U, V \in \mathbb{R}^{k \times d}$ and let \mathcal{N}_d be the standard Gaussian in d dimensions. We have that

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\sigma_1(\boldsymbol{U}\boldsymbol{x}) \neq \sigma_1(\boldsymbol{V}\boldsymbol{x})] \leq c \cdot k \cdot \sqrt{\log k} \cdot \max_{i \neq j} \theta(\boldsymbol{U}_i - \boldsymbol{U}_j, \boldsymbol{V}_i - \boldsymbol{V}_j) \,,$$

- where c > 0 is some universal constant.
- ⁸⁹¹ *Proof.* We have that

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\sigma_1(\boldsymbol{U}\boldsymbol{x}) \neq \sigma_1(\boldsymbol{V}\boldsymbol{x})] = \sum_{i \in [k]} \Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\sigma_1(\boldsymbol{U}\boldsymbol{x}) = i, \sigma_1(\boldsymbol{V}\boldsymbol{x}) \neq i] \,.$$

We have that $C_{U}^{(i)} = \mathbf{1}\{\mathbf{x} : \sigma_1(\mathbf{U}\mathbf{x}) = i\} = \prod_{j \neq i} \mathbf{1}\{(\mathbf{U}_i - \mathbf{U}_j) \cdot \mathbf{x} \ge 0\}$ is the set indicator of a homogeneous polyhedral cone as the intersection of k - 1 homogeneous halfspaces. Similarly, we consider the cone $C_{V}^{(i)} = \{\mathbf{x} : \sigma_1(\mathbf{V}\mathbf{x}) = i\}$. Hence, we have that $\{\mathbf{x} : \sigma_1(\mathbf{V}\mathbf{x}) \neq i\}$ is the complement of a homogeneous polyhedral cone. Let us define $C_{U}^{(i)} : \mathbb{R}^d \mapsto \{0, 1\}$ and $C_{V}^{(i)} : \mathbb{R}^d \mapsto \{0, 1\}$ be the associated indicator functions of the two cones. We have that

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[\sigma_1(\boldsymbol{U}\boldsymbol{x})=i,\sigma_1(\boldsymbol{V}\boldsymbol{x})\neq i] = \Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[C_{\boldsymbol{U}}^{(i)}(\boldsymbol{x})=1,C_{\boldsymbol{V}}^{(i)}(\boldsymbol{x})=0].$$

⁸⁹⁷ Finally, we have that

$$\mathcal{C}_{\boldsymbol{U}}^{(i)} \cap \left(\mathcal{C}_{\boldsymbol{V}}^{(i)}\right)^{c} = \mathcal{C}_{\boldsymbol{U}}^{(i)} \setminus \mathcal{C}_{\boldsymbol{V}}^{(i)} \subseteq \mathcal{C}_{\boldsymbol{U}}^{(i)} \setminus \mathcal{C}_{\boldsymbol{V}}^{(i)} \cup \mathcal{C}_{\boldsymbol{V}}^{(i)} \setminus \mathcal{C}_{\boldsymbol{U}}^{(i)} .$$

⁸⁹⁸ We can hence apply Lemma 13 for the cones $C_U^{(i)}, C_V^{(i)}$ for each $i \in [k]$.

Lemma 13 (Cone Disagreement). Let
$$C_1 : \mathbb{R}^d \mapsto \{0, 1\}$$
 be the indicator function of the homoge-
neous polyhedral cone defined by the k unit vectors $v_1, \ldots, v_k \in \mathbb{R}^d$, i.e., $C_1(x) = \prod_{i=1}^k \mathbf{1}\{v_i \cdot x \ge 0\}$.
Similarly, define $C_2 : \mathbb{R}^d \mapsto \{0, 1\}$ to be the homogeneous polyhedral cone with normal vectors

902 u_1, \ldots, u_k . It holds that

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d} [C_1(\boldsymbol{x}) \neq C_2(\boldsymbol{x})] \le c \sqrt{\log(k)} \, \max_{i \in [k]} \theta(\boldsymbol{v}_i, \boldsymbol{u}_i) \,,$$

903 where c > 0 is some universal constant.

Proof. To simplify notation, denote $\theta = \max_{i \in [k]} \theta(v_i, u_i)$. We first observe that it suffices to prove 904 the upper bound on the probability of $C_1(\mathbf{x}) \neq C_2(\mathbf{x})$ for sufficiently small values of θ . Indeed, if 905 we have that the bound is true for θ smaller than some θ_0 we can then form a path of sufficiently large 906 length N (in particular we need $\theta/N \le \theta_0$) starting from the vectors v_1, \ldots, v_k to the final vectors 907 u_1, \ldots, u_k , where at each step we only rotate the vectors by at most $\theta/N \leq \theta_0$. By the triangle 908 inequality, we immediately obtain that the probability that $C_1(x) \neq C_2(x)$ is at most equal to the sum of the probabilities of the intermediate steps which is at most $\sum_{i=1}^{N} c\sqrt{\log(k)} \frac{\theta}{N} = c\sqrt{\log(k)}\theta$. 909 910 Notice in the above argument the constant θ_0 can be arbitrarily small and may also depend on k and 911 912 d.

We define the indicator of the positive orthant in k dimensions to be $R(t) = \prod_{i=1}^{k} \mathbf{1}\{t_i \ge 0\}$. Using this notation, we have that the cone indicator can be written as $C_1(x) = R(v_1 \cdot x, \dots, v_k \cdot x) =$ R(Vx), where V is the $k \times d$ matrix whose *i*-th row is the vector v_i . Moreover, we define the *i*-th face of the cone R(Vx) to be

$$F_i(\boldsymbol{V} \boldsymbol{x}) = R(\boldsymbol{V} \boldsymbol{x}) \mathbf{1} \{ \boldsymbol{v}_i \cdot \boldsymbol{x} = 0 \}$$
.

We will first handle the case where only one of the normal vectors v_i changes. We show the following claim.

919 **Claim 10.** Let $v_1, \ldots, v_k \in \mathbb{R}^d$ and $r \in \mathbb{R}^d$ with $\theta(v_1, r) \leq \theta$ for some sufficiently small $\theta \in (0, \pi/2)$. It holds that

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}\left[R(\boldsymbol{v}_1\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})\neq R(\boldsymbol{r}\cdot\boldsymbol{x},\boldsymbol{v}_2\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})\right]\leq c\cdot\theta\cdot\Gamma(F_1)\sqrt{\log\left(\frac{1}{\Gamma(F_1)}+1\right)}$$

where F_1 is the face with $v_1 \cdot x = 0$ of the cone R(Vx) and c is some universal constant.



Figure 1: The vectors r, v_1 and q and the disagreement region of the halfspaces with normal vectors r and v_1 .

922 *Proof.* We have

$$\begin{split} & \Pr_{\boldsymbol{x} \sim \mathcal{N}_d} \left[R(\boldsymbol{v}_1 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \neq R(\boldsymbol{r} \cdot \boldsymbol{x}, \boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \right] \\ &= \mathop{\mathbf{E}}_{\boldsymbol{x} \sim \mathcal{N}_d} \left[\left| R(\boldsymbol{v}_1 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) - R(\boldsymbol{r} \cdot \boldsymbol{x}, \boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \right| \right] \\ &= \mathop{\mathbf{E}}_{\boldsymbol{x} \sim \mathcal{N}_d} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \left| \mathbf{1} \{ \boldsymbol{v}_1 \cdot \boldsymbol{x} \geq 0 \} - \mathbf{1} \{ \boldsymbol{r} \cdot \boldsymbol{x} \geq 0 \} \right| \right]. \end{split}$$

We have that $|\mathbf{1}\{v_1 \cdot x \ge 0\} - \mathbf{1}\{r \cdot x \ge 0\}| = \mathbf{1}\{(v_1 \cdot x)(r \cdot x) < 0\}$, i.e., this is the event that the halfspaces $\mathbf{1}\{v_1 \cdot x \ge 0\}$ and $\mathbf{1}\{r \cdot x \ge 0\}$ disagree. Let q be the normalized projection of r onto the orthogonal complement of v_1 , i.e., $q = \operatorname{proj}_{v_1^{\perp}} r/||\operatorname{proj}_{v_1^{\perp}} r||_2$. We have that v_1 and q is an orthonormal basis of the subspace spanned by the vectors v_1 and r. We have that $r = \cos \theta(v_1, r)v_1 + \sin \theta(v_1, r)q$. Moreover, we have that the region $(v_1 \cdot x)(r \cdot x) < 0$ is equal to

$$\left\{0 < \boldsymbol{v}_1 \cdot \boldsymbol{x} < -(\boldsymbol{q} \cdot \boldsymbol{x}) \tan \theta(\boldsymbol{v}_1, \boldsymbol{r})\right\} \cup \left\{-(\boldsymbol{q} \cdot \boldsymbol{x}) \tan \theta(\boldsymbol{v}_1, \boldsymbol{r}) < \boldsymbol{v}_1 \cdot \boldsymbol{x} < 0\right\}.$$

Thus, we have that the disagreement region $(v_1 \cdot x)(r \cdot x) < 0$ is a subset of the region $\{|v_1 \cdot x| \le |q \cdot x| \tan \theta(v_1, r)\}$. Since $\tan \theta(v_1, r) \le \theta$ and we have that θ is sufficiently small we can also replace the above region by the larger region: $\{|v_1 \cdot x| \le 2\theta |q \cdot x|\}$. Therefore, we have

$$\begin{split} & \underset{\boldsymbol{x} \sim \mathcal{N}_d}{\mathbf{E}} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \ \mathbf{1}\{(\boldsymbol{v}_1 \cdot \boldsymbol{x})(\boldsymbol{r} \cdot \boldsymbol{x}) < 0\}\} \right] \\ & \qquad \leq \underset{\boldsymbol{x} \sim \mathcal{N}_d}{\mathbf{E}} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \ \mathbf{1}\{|\boldsymbol{v}_1 \cdot \boldsymbol{x}| \leq 2\theta | \boldsymbol{q} \cdot \boldsymbol{x}|\} \right] \,. \end{split}$$

⁹³² The derivative of the above expression with respect to θ is equal to

$$\mathbf{E}_{\boldsymbol{x}\sim\mathcal{N}_d}\left[R(\boldsymbol{v}_2\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})\;\delta\left(\frac{|\boldsymbol{v}_1\cdot\boldsymbol{x}|}{2|\boldsymbol{q}\cdot\boldsymbol{x}|}-\theta\right)\right],$$

where $\delta(t)$ is the Dirac delta function. At $\theta = 0$ and using the property that $\delta(t/a) = a\delta(t)$, we have that the above derivative is equal to

$$2 \mathop{\mathbf{E}}_{\boldsymbol{x} \sim \mathcal{N}_d} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \mid \boldsymbol{q} \cdot \boldsymbol{x} \mid \delta(|\boldsymbol{v}_1 \cdot \boldsymbol{x}|) \right] \,.$$

- Notice that, if we did not have the term $|\mathbf{q} \cdot \mathbf{x}|$, the above expression would be exactly equal to two times the Gaussian surface area of the face with $\mathbf{v}_1 \cdot \mathbf{x} = 0$, i.e., it would be equal to $2\Gamma(F_1)$. We
- times the Gaussian surface area of the face with $v_1 \cdot x = 0$, i.e., it would be equal to $2\Gamma(F_1)$. We now show that this extra term of $|q \cdot x|$ can only increase the above surface integral by at most a

938 logarithmic factor. We have that

$$\begin{split} & \underset{\boldsymbol{x} \sim \mathcal{N}_d}{\mathbf{E}} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \left| \boldsymbol{q} \cdot \boldsymbol{x} \right| \delta(|\boldsymbol{v}_1 \cdot \boldsymbol{x}|) \right] = \int_{\boldsymbol{x} \in F_1} \phi_d(\boldsymbol{x}) |\boldsymbol{q} \cdot \boldsymbol{x}| d\mu(\boldsymbol{x}) \\ & \leq \int_{\boldsymbol{x} \in F_1} \phi_d(\boldsymbol{x}) |\boldsymbol{q} \cdot \boldsymbol{x}| \mathbf{1} \{ |\boldsymbol{q} \cdot \boldsymbol{x}| \leq \xi \} d\mu(\boldsymbol{x}) + \int_{\boldsymbol{x} \in F_1} \phi_d(\boldsymbol{x}) |\boldsymbol{q} \cdot \boldsymbol{x}| \mathbf{1} \{ |\boldsymbol{q} \cdot \boldsymbol{x}| \geq \xi \} d\mu(\boldsymbol{x}) \\ & \leq \xi \int_{\boldsymbol{x} \in F_1} \phi_d(\boldsymbol{x}) d\mu(\boldsymbol{x}) + \int_{\boldsymbol{x} \in F_1} \phi_d(\boldsymbol{x}) |\boldsymbol{q} \cdot \boldsymbol{x}| \mathbf{1} \{ |\boldsymbol{q} \cdot \boldsymbol{x}| \geq \xi \} d\mu(\boldsymbol{x}) \,, \end{split}$$

where $d\mu(\mathbf{x})$ is the standard surface measure in \mathbb{R}^d . The first term above is exactly equal to the Gaussian surface area of the face F_1 . To bound from above the second term we can use the fact that the face F_1 is a subset of the hyperplane $\mathbf{v}_1 \cdot \mathbf{x} = 0$, i.e., it holds that $F_1 \subseteq {\mathbf{x} : |\mathbf{v}_1 \cdot \mathbf{x}| = 0}$. To simplify notation we may assume that $\mathbf{v}_1 = \mathbf{e}_1$ and $\mathbf{q} = \mathbf{e}_2$ (recall that \mathbf{v}_1 and \mathbf{q} are orthogonal unit vectors), and in this case we obtain

$$\begin{split} \int_{\boldsymbol{x}\in F_1} \phi_d(\boldsymbol{x}) |\boldsymbol{q}\cdot\boldsymbol{x}| \mathbf{1}\{|\boldsymbol{q}\cdot\boldsymbol{x}| \geq \xi\} d\mu(\boldsymbol{x}) &\leq \int_{x_1=0} \phi_d(\boldsymbol{x}) |x_2| \mathbf{1}\{|x_2| \geq \xi\} d\mu(\boldsymbol{x}) \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} |x_2| \mathbf{1}\{|x_2| \geq \xi\} \frac{e^{-x_2^2/2}}{\sqrt{2\pi}} dx_2 \\ &= \frac{1}{\pi} e^{-\xi^2/2} \,. \end{split}$$

Combining the above bounds we obtain that the derivative with respect to θ of the expression $\mathbf{E}_{\boldsymbol{x}\sim\mathcal{N}_d} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \mathbf{1} \{ |\boldsymbol{v}_1 \cdot \boldsymbol{x}| \leq 2\theta | \boldsymbol{q} \cdot \boldsymbol{x}| \} \right]$ is equal to

$$\frac{d}{d\theta} \Big(\mathbf{E}_{\boldsymbol{x} \sim \mathcal{N}_d} \left[R(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \, \mathbf{1} \{ |\boldsymbol{v}_1 \cdot \boldsymbol{x}| \le 2\theta | \boldsymbol{q} \cdot \boldsymbol{x}| \} \right] \Big) \Big|_{\theta=0} \le 2\xi \Gamma(F_1) + \frac{2e^{-\xi^2/2}}{\pi}$$

By picking $\xi = \sqrt{2 \log(1 + 1/\Gamma(F_1))}$, the result follows since up to introducing $o(\theta)$ error we can bound the term $\Pr_{\boldsymbol{x} \sim \mathcal{N}_d} [R(\boldsymbol{v}_1 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \neq R(\boldsymbol{r} \cdot \boldsymbol{x}, \boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x})]$ by its derivative with respect to θ (evaluated at 0) times θ .

We can complete the proof of Lemma 13 using Claim 10. In order to bound the disagreement of the cones C_1 and C_2 we can start from C_1 and change one of its vectors at a time so that we can use Claim 10 that can handle this case. For example, at the first step, we can swap v_1 for u_1 and use the triangle inequality to obtain that

$$\begin{split} \mathbf{Pr}_{\boldsymbol{x}\sim\mathcal{N}_d}[C_1(\boldsymbol{x})\neq C_2(\boldsymbol{x})] &\leq \mathbf{Pr}_{\boldsymbol{x}\sim\mathcal{N}_d}[R(\boldsymbol{v}_1\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})\neq R(\boldsymbol{u}_1\cdot\boldsymbol{x},\boldsymbol{v}_2\cdot\boldsymbol{x}\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})] \\ &+ \mathbf{Pr}_{\boldsymbol{x}\sim\mathcal{N}_d}[R(\boldsymbol{u}_1\cdot\boldsymbol{x},\boldsymbol{v}_2\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})\neq R(\boldsymbol{u}_1\cdot\boldsymbol{x},\boldsymbol{u}_2\cdot\boldsymbol{x}\ldots,\boldsymbol{u}_k\cdot\boldsymbol{x})] \\ &\leq c\cdot\theta\;\Gamma(F_1)\sqrt{\log(1/\Gamma(F_1)+1)} \\ &+ \mathbf{Pr}_{\boldsymbol{x}\sim\mathcal{N}_d}[R(\boldsymbol{u}_1\cdot\boldsymbol{x},\boldsymbol{v}_2\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})\neq R(\boldsymbol{u}_1\cdot\boldsymbol{x},\boldsymbol{u}_2\cdot\boldsymbol{x}\ldots,\boldsymbol{u}_k\cdot\boldsymbol{x})]\,, \end{split}$$

where $F_1 = F_1(Vx)$ is the face with $v_1 \cdot x = 0$ of the cone C_1 . Notice that we have replaced v_1 by u_1 in the above bound. Our plan is to use the triangle inequality and continue replacing the vectors of C_1 by the vectors of C_2 sequentially. To make this formal we define the matrix $A^{(i)} \in \mathbb{R}^{k \times d}$ whose first i - 1 rows are the vectors u_1, \ldots, u_{i-1} and its last k - i + 1 rows are the vectors v_i, \ldots, v_k , i.e.,

$$oldsymbol{A}_{j}^{(i)} = egin{cases} oldsymbol{u}_{j} & ext{if} \quad 1 \leq j \leq i-1, \ oldsymbol{v}_{j} & ext{if} \quad i \leq j \leq k \,. \end{cases}$$

Notice that $A^{(1)} = V$ and $A^{(k+1)} = U$. Using the triangle inequality we obtain that

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[C_1(\boldsymbol{x})\neq C_2(\boldsymbol{x})]\leq \sum_{i=1}^k\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[R(\boldsymbol{A}^{(i)}\boldsymbol{x})\neq R(\boldsymbol{A}^{(i+1)}\boldsymbol{x})].$$

Since the matrices $A^{(i)}$ and $A^{(i+1)}$ only differ on one row, we can use Claim 10 to obtain the following bound:

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[C_1(\boldsymbol{x}) \neq C_2(\boldsymbol{x})] \le c \cdot \theta \cdot \sum_{i=1}^k \Gamma(F_i(\boldsymbol{A}^{(i)}\boldsymbol{x})) \sqrt{\frac{1}{\Gamma(F_i(\boldsymbol{A}^{(i)}\boldsymbol{x}))} + 1}.$$

We now observe that the Gaussian surface area $\Gamma(F_i(\mathbf{A}^{(i)}\mathbf{x}))$ is a continuous function of the matrix $\mathbf{A}^{(i)}$. By flattening the matrix $\mathbf{A}^{(i)}$ (since it is isomorphic to a vector $\mathbf{z} \in \mathbb{R}^{n^2}$) and letting $S_{\mathbf{z}}$ be the induced surface { $\mathbf{x} : R(\mathbf{A}^{(i)}\mathbf{x}) = 1 \land \mathbf{v}_i \cdot \mathbf{x} = 0$ }, it suffices to show that

$$\lim_{\boldsymbol{w}\to\boldsymbol{z}}\int\phi_n(\boldsymbol{x})\mathbf{1}\{\boldsymbol{x}\in S_{\boldsymbol{w}}\}d\mu(\boldsymbol{x})=\int\phi_n(\boldsymbol{x})\mathbf{1}\{\boldsymbol{x}\in S_{\boldsymbol{z}}\}d\mu(\boldsymbol{x})\,,$$

by the smoothness of the surface S_z . Consider a sequence of functions (g_m) and vectors (w_m) so that $g_m(x) = \phi_n(x) \mathbf{1}\{x \in S_{w_m}\}$ and $\lim_{m\to\infty} w_m = z$. Note that $|g_m(x)| \leq 1$ everywhere. Hence, by the dominated convergence theorem, we have that

$$\lim_{m\to\infty}\int g_m(\boldsymbol{x})d\mu(\boldsymbol{x})=\int\lim_{m\to\infty}g_m(\boldsymbol{x})d\mu(\boldsymbol{x})=\int\phi_n(\boldsymbol{x})\lim_{m\to\infty}\mathbf{1}\{\boldsymbol{x}\in S_{\boldsymbol{w}_m}\}d\mu(\boldsymbol{x})\,.$$

Since the sequence consists of smooth surfaces, we have that $\lim_{m\to\infty} \mathbf{1}\{x \in S_{w_m}\} = \mathbf{1}\{x \in S_z\}$ and so the Gaussian surface area is continuous with respect to the matrix $A^{(i)}$ for any $i \in [k]$.

Also, as $\theta \to 0$, we have that $A^{(i)} \to V$. This is because the sequence of matrices $A^{(i)}$ depends only on the vectors u_j and v_j for $j \in [k]$ and the following two properties hold true: $\theta = \max_{j \in [k]} \theta(v_j, u_j)$ and all the vectors are unit. Hence, as θ tends to zero, they tend to become the same vectors and so any matrix $A^{(i)}$ tends to become V. Therefore, taking this limit we obtain that for $\theta \to 0$ it holds that

$$\lim_{k \to 0} \frac{\mathbf{Pr}_{\boldsymbol{x} \sim \mathcal{N}_d}[C_1(\boldsymbol{x}) \neq C_2(\boldsymbol{x})]}{\theta} \le c \cdot \sum_{i=1}^k \Gamma(F_i(\boldsymbol{V}\boldsymbol{x})) \sqrt{\log\left(1/\Gamma(F_i(\boldsymbol{V}\boldsymbol{x})) + 1\right)}.$$
(1)

We will now use the following lemma that shows that the surface area of any homogeneous polyhedral cone is independent of the number of faces k and in fact is at most 1 for all k.

Lemma 14 (Gaussian Surface Area of Homogeneous Cones [Naz03]). Let C be a cone with apex at the origin (i.e., an intersection of arbitrarily many halfspaces all of whose boundaries contain the origin). Then C has Gaussian surface area $\Gamma(C)$ at most 1.

979 Using Lemma 14 we obtain that $\sum_{i=1}^{k} \Gamma(F_i(\mathbf{Vx})) \leq 1$. Next, we observe that, when 980 the positive numbers a_1, \ldots, a_k satisfy $\sum_{i=1}^{k} a_i \leq 1$, it holds that $\sum_{i=1}^{k} a_i \sqrt{\log(1/a_i)} \leq 1$

981 $\sqrt{\sum_{i=1}^{k} a_i \log(1/a_i)} \le \sqrt{\log(k)}$ (using the fact that the uniform distribution maximizes the entropy). 982 Using this fact and Equation (1), we obtain

$$\lim_{\theta \to 0} \frac{\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[C_1(\boldsymbol{x}) \neq C_2(\boldsymbol{x})]}{\theta} \le c\sqrt{\log(k)} \,.$$

Thus, we have shown that, for sufficiently small θ , it holds that $\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[C_1(\boldsymbol{x}) \neq C_2(\boldsymbol{x})] \leq c\sqrt{\log(k)}\theta$, but, as we discussed in the start of the proof, the general bound follows directly from the bound for sufficiently small values of $\theta > 0$.

⁹⁸⁶ C Learning in Top-*r* Disagreement from Label Rankings

1

We prove the next result which corresponds to a proper learning algorithm for LSF in the presence of bounded noise with respect to the top-*r* disagreement.

Theorem 6 (Proper Top-*r* Learning Algorithm). Fix $\eta \in [0, 1/2)$, $r \in [k]$ and $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2. There exists an algorithm that draws $N = \widetilde{O}\left(\frac{d \ rk}{\epsilon(1-2\eta)^6}\log(1/\delta)\right)$ samples from \mathcal{D} , runs in poly(N) time and, with probability at least $1 - \delta$, outputs a Linear Sorting function $h : \mathbb{R}^d \to \mathbb{S}_k$ that is ϵ -close in top-rdisagreement to the target.

- The main result of this section is the next lemma, which directly implies the above theorem (using the same steps as the proof of Theorem 5).
- **Lemma 15** (Top-*r* Misclassification). Let $r \in [k]$. Consider two matrices $U, V \in \mathbb{R}^{k \times d}$ and let \mathcal{N}_d be the standard Gaussian in d dimensions. We have that

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[\sigma_{1..r}(\boldsymbol{U}\boldsymbol{x}) \neq \sigma_{1..r}(\boldsymbol{V}\boldsymbol{x})] \le c \cdot k \cdot r \cdot \sqrt{\log(kr)} \cdot \max_{i \neq j} \theta(\boldsymbol{U}_i - \boldsymbol{U}_j, \boldsymbol{V}_i - \boldsymbol{V}_j),$$

998 where c > 0 is some universal constant.

Proof. Let us set $\sigma_{1..r}(Wx)$ denote the ordering of the top-*r* alternatives in the ranking $\sigma(Wx)$. Moreover, recall that $\sigma_{\ell}(Wx)$ denotes the alternative in the ℓ -th position of the ranking $\sigma(Wx)$. For two matrices $U, V \in \mathbb{R}^{k \times d}$, we have that

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[\sigma_{1..r}(\boldsymbol{U}\boldsymbol{x})\neq\sigma_{1..r}(\boldsymbol{V}\boldsymbol{x})]=\sum_{j=1}^k\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}\left[\bigcup_{\ell=1}^r\{j=\sigma_\ell(\boldsymbol{U}\boldsymbol{x}),j\neq\sigma_\ell(\boldsymbol{V}\boldsymbol{x})\}\right].$$

The first step is to understand the geometry of the set $\bigcup_{\ell=1}^{r} \{ \boldsymbol{x} : j = \sigma_{\ell}(\boldsymbol{U}\boldsymbol{x}) \} = \{ \boldsymbol{x} : j \in \sigma_{1..r}(\boldsymbol{U}\boldsymbol{x}) \}$ for $j \in [k]$. We have that this set is equal to

$$\mathcal{T}_{\boldsymbol{U}}^{(j)} = \bigcup_{S \subseteq [k]: |S| \le r-1} \bigcap_{i \in S} \{ \boldsymbol{x} : (\boldsymbol{U}_i - \boldsymbol{U}_j) \cdot \boldsymbol{x} \ge 0 \} \cap \bigcap_{i \notin S} \{ \boldsymbol{x} : (\boldsymbol{U}_i - \boldsymbol{U}_j) \cdot \boldsymbol{x} \le 0 \}$$

In words, $\mathcal{T}_{U}^{(j)}$ iterates over any possible collection of alternatives that can win the element j(they lie in the set of top elements S) and the remaining elements lose when compared with j(they lie in the complement set $[k] \setminus S$). Overloading the notation, let us define the mapping $T(t) = T(t_1, ..., t_k) = \sum_{S \subseteq [k]: |S| \le r-1} \prod_{i \in S} \mathbf{1}\{t_i \ge 0\} \prod_{i \notin S} \mathbf{1}\{t_i \le 0\}$. Using this mapping, we can define the indicator of the set $T_{U}^{(j)}$ as $T((U_1 - U_j) \cdot \boldsymbol{x}, ..., (U_k - U_j) \cdot \boldsymbol{x})$. The top-rdisagreement $\mathbf{Pr}_{\boldsymbol{x} \sim \mathcal{N}_d}[j \in \sigma_{1..r}(\boldsymbol{U}\boldsymbol{x}), j \notin \sigma_{1..r}(\boldsymbol{V}\boldsymbol{x})]$ is equal to:

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[T((\boldsymbol{U}_1-\boldsymbol{U}_j)\cdot\boldsymbol{x},...,(\boldsymbol{U}_k-\boldsymbol{U}_j)\cdot\boldsymbol{x})=1,T((\boldsymbol{V}_1-\boldsymbol{V}_j)\cdot\boldsymbol{x},...,(\boldsymbol{V}_k-\boldsymbol{V}_j)\cdot\boldsymbol{x})=0].$$

1010 So we have that

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[\sigma_{1..r}(\boldsymbol{U}\boldsymbol{x})\neq\sigma_{1..r}(\boldsymbol{V}\boldsymbol{x})] = \sum_{j=1}^k \Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[T_j(\boldsymbol{U}\boldsymbol{x})=1, T_j(\boldsymbol{V}\boldsymbol{x})=0] \le \sum_{j=1}^k \Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[T_j(\boldsymbol{U}\boldsymbol{x})\neq T_j(\boldsymbol{V}\boldsymbol{x})].$$

1011 In order to show the desired bound, it suffices to prove the following two lemmas.

1012 **Lemma 16** (Disagreement Region). Consider a positive integer $r \le k$. For any $j \in [k]$, it holds that

$$\lim_{\theta \to 0} \frac{\mathbf{Pr}_{\boldsymbol{x} \sim \mathcal{N}_d}[T_j(\boldsymbol{U}\boldsymbol{x}) \neq T_j(\boldsymbol{V}\boldsymbol{x})]}{\theta} \le c \cdot \sum_{i \in [k]} \Gamma(F_i^j) \sqrt{\log\left(\frac{1}{\Gamma(F_i^j)} + 1\right)},$$

where c > 0 is some constant and F_i^j is the surface $\{ \boldsymbol{x} : j \in \sigma_{1..r}(\boldsymbol{V}\boldsymbol{x}) \} \cap \{ \boldsymbol{x} : \boldsymbol{V}_i \cdot \boldsymbol{x} = \boldsymbol{V}_j \cdot \boldsymbol{x} \}$ for the matrix $\boldsymbol{V} \in \mathbb{R}^{k \times d}$.

1015 and,

1016 **Lemma 17.** Let F_i^j , r, k as in the previous lemma. It holds that

$$\sum_{i \in [k]} \sum_{j \in [k]} \Gamma(F_i^j) \le 2kr.$$

1017 Applying these two lemmas, we get that

$$Z := \lim_{\theta \to 0} \frac{\sum_{j \in [k]} \mathbf{Pr}_{\boldsymbol{x} \sim \mathcal{N}_d}[T_j(\boldsymbol{U}\boldsymbol{x}) \neq T_j(\boldsymbol{V}\boldsymbol{x})]}{\theta} \le c \cdot \sum_{j \in [k]} \sum_{i \in [k]} \Gamma(F_i^j) \sqrt{\log\left(\frac{1}{\Gamma(F_i^j)} + 1\right)}.$$

Let us set $\Gamma'(F_i^j) = \Gamma(F_i^j)/(2kr)$. Then we have that 1018

$$Z \leq 2ckr \cdot \sum_{j \in [k]} \sum_{i \in [k]} \Gamma'(F_i^j) \sqrt{\log\left(\frac{1}{2kr \cdot \Gamma'(F_i^j)} + 1\right)}.$$

It suffices to bound the quantity 1019

$$\sum_{j \in [k]} \sum_{i \in [k]} \Gamma'(F_i^j) \sqrt{\log\left(\frac{1}{\Gamma'(F_i^j)} + 1\right)} = O\left(kr\sqrt{\log(kr)}\right) \,,$$

where we used a similar "entropy-like" inequality as we did in the top-1 case. This yields (by recalling 1020 that it is sufficient to consider only the case of arbitrarily small angles, as in the top-1 case) that 1021

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d} [\sigma_{1..r}(\boldsymbol{U}\boldsymbol{x}) \neq \sigma_{1..r}(\boldsymbol{V}\boldsymbol{x})] \le c \ rk \ \sqrt{\log(kr)} \cdot \max_{i \ne j} \theta(\boldsymbol{U}_i - \boldsymbol{U}_j, \boldsymbol{V}_i - \boldsymbol{V}_j) ,$$

for some universal constant c. 1022

C.1 The proof of Lemma 16 1023

We proceed with the proof of the key lemma concerning the disagreement region. We first show the 1024 following claim where we only change a single vector. Recall that 1025

$$T(\boldsymbol{V}\boldsymbol{x}) = \sum_{S:|S| \le r-1} \prod_{i \in S} \mathbf{1}\{\boldsymbol{v}_i \cdot \boldsymbol{x} \ge 0\} \prod_{i \notin S} \mathbf{1}\{\boldsymbol{v}_i \cdot \boldsymbol{x} \le 0\}.$$

- We will be interested in the surface $F_1 := F_1(Vx) = T(Vx)\mathbf{1}\{v_1 \cdot x = 0\}.$ 1026
- Claim 11. Let $v_1, \ldots, v_k \in \mathbb{R}^d$ and $r \in \mathbb{R}^d$ with $\theta(v_1, r) \leq \theta$ for some sufficiently small $\theta \in (0, \pi/2)$. It holds that 1027 1028

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[T(\boldsymbol{v}_1\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})\neq T(\boldsymbol{r}\cdot\boldsymbol{x},\boldsymbol{v}_2\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_k\cdot\boldsymbol{x})]\leq c\cdot\theta\cdot\Gamma(F_1)\sqrt{\log\left(\frac{1}{\Gamma(F_1)}+1\right)},$$

where F_1 is the surface $T(\mathbf{V}\mathbf{x}) \cap \{\mathbf{x} : \mathbf{v}_1 \cdot \mathbf{x} = 0\}$ and c is some universal constant. 1029

Proof. We first decompose the sum of T(Vx) depending on whether $1 \in S$ or not. Hence, we have 1030 that $T(v_1 \cdot x, \ldots, v_k \cdot x) = T^+(v_1 \cdot x, \ldots, v_k \cdot x) + T^-(v_1 \cdot x, \ldots, v_k \cdot x)$ where 1031

$$T^{+}(\boldsymbol{v}_{1}\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_{k}\cdot\boldsymbol{x}) = \sum_{S\subseteq[k]:|S|\leq r-1,1\in S}\prod_{i\in S}\mathbf{1}\{\boldsymbol{v}_{i}\cdot\boldsymbol{x}\geq 0\}\prod_{i\notin S}\mathbf{1}\{\boldsymbol{v}_{i}\cdot\boldsymbol{x}\leq 0\}$$
$$= \sum_{S\subseteq[k]:|S|\leq r-1,1\in S}\mathbf{1}\{\boldsymbol{v}_{1}\cdot\boldsymbol{x}\geq 0\}\cdot\prod_{i\in S\setminus\{1\}}\mathbf{1}\{\boldsymbol{v}_{i}\cdot\boldsymbol{x}\geq 0\}\prod_{i\notin S}\mathbf{1}\{\boldsymbol{v}_{i}\cdot\boldsymbol{x}\leq 0\}$$
$$= \mathbf{1}\{\boldsymbol{v}_{1}\cdot\boldsymbol{x}\geq 0\}\cdot\sum_{S\subseteq[k]:|S|\leq r-1,1\in S}\prod_{i\in S\setminus\{1\}}\mathbf{1}\{\boldsymbol{v}_{i}\cdot\boldsymbol{x}\geq 0\}\prod_{i\notin S}\mathbf{1}\{\boldsymbol{v}_{i}\cdot\boldsymbol{x}\leq 0\}$$
$$=:\mathbf{1}\{\boldsymbol{v}_{1}\cdot\boldsymbol{x}\geq 0\}\cdot G^{+}(\boldsymbol{v}_{2}\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_{k}\cdot\boldsymbol{x}),$$

and similarly 1032

$$T^-(oldsymbol{v}_1 \cdot oldsymbol{x}, \dots, oldsymbol{v}_k \cdot oldsymbol{x}) = \mathbf{1}\{oldsymbol{v}_1 \cdot oldsymbol{x} \leq 0\} \cdot \sum_{S \subseteq [k]: |S| \leq r-1, 1 \notin S} \prod_{i \in S} \mathbf{1}\{oldsymbol{v}_i \cdot oldsymbol{x} \geq 0\} \prod_{i \notin S \setminus \{1\}} \mathbf{1}\{oldsymbol{v}_i \cdot oldsymbol{x} \leq 0\}$$

=: $\mathbf{1}\{oldsymbol{v}_1 \cdot oldsymbol{x} \leq 0\} \cdot G^-(oldsymbol{v}_2 \cdot oldsymbol{x}, \dots, oldsymbol{v}_k \cdot oldsymbol{x}).$

Notice that the indicator G^s does not depend on the alternative 1 for $s \in \{-,+\}$. Since $T : \mathbb{R}^k \to \mathbb{R}^k$ 1033 $\{0,1\}$, we have that 1034

$$\begin{aligned} & \Pr_{\boldsymbol{x} \sim \mathcal{N}_d} [T(\boldsymbol{v}_1 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \neq T(\boldsymbol{r} \cdot \boldsymbol{x}, \boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x})] \\ &= \mathop{\mathbf{E}}_{\boldsymbol{x} \sim \mathcal{N}_d} [|T(\boldsymbol{v}_1 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) - T(\boldsymbol{r} \cdot \boldsymbol{x}, \boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x})|] \\ &\leq \sum_{s \in \{-,+\}} \mathop{\mathbf{E}}_{\boldsymbol{x} \sim \mathcal{N}_d} [|T^s(\boldsymbol{v}_1 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) - T^s(\boldsymbol{r} \cdot \boldsymbol{x}, \boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x})|] \\ &= \sum_{s \in \{-,+\}} \mathop{\mathbf{E}}_{\boldsymbol{x} \sim \mathcal{N}_d} [G^s(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \cdot |\mathbf{1}\{s \cdot \boldsymbol{v}_1 \cdot \boldsymbol{x} \ge 0\} - \mathbf{1}\{s \cdot \boldsymbol{r} \cdot \boldsymbol{x} \ge 0\}|] \end{aligned}$$

Let us focus on the case s = +. The difference between the two indicators in the last line of the 1035 above equation corresponds to the event that the halfspaces $\mathbf{1}\{v_1 \cdot x > 0\}$ and $\mathbf{1}\{r \cdot x > 0\}$ disagree. 1036 Hence, we have that $|\mathbf{1}\{v_1 \cdot x \ge 0\} - \mathbf{1}\{r \cdot x \ge 0\}| = \mathbf{1}\{(v_1 \cdot x)(r \cdot x) < 0\}$. Note that the above 1037 indicator depends on both v_1 and r. We would like to work only with one of these two vectors. To 1038 this end, let us introduce q, the normalized projection of r onto the orthogonal complement of v_1 , i.e., 1039 $q = \text{proj}_{v_1^{\perp}} r / \|\text{proj}_{v_1^{\perp}} r\|_2$. We have that v_1 and q is an orthonormal basis of the subspace spanned 1040 by the vectors v_1 and r. Notice that $r = \cos \theta(v_1, r)v_1 + \sin \theta(v_1, r)q$, by the construction of q. 1041 Our goal is to understand the structure of the region $(v_1 \cdot x)(r \cdot x) < 0$. This set is equal to 1042

$$\{0 < \boldsymbol{v}_1 \cdot \boldsymbol{x} < -(\boldsymbol{q} \cdot \boldsymbol{x}) \tan \theta(\boldsymbol{v}_1, \boldsymbol{r})\} \cup \{-(\boldsymbol{q} \cdot \boldsymbol{x}) \tan \theta(\boldsymbol{v}_1, \boldsymbol{r}) < \boldsymbol{v}_1 \cdot \boldsymbol{x} < 0\}.$$

To see this, we have that $(v_1 \cdot x)(r \cdot x) = (v_1 \cdot x)(\cos \theta(v_1, r)v_1 \cdot x + \sin \theta(v_1, r)q \cdot x)$. This quantity must be negative. The left-hand set considers the case where $v_1 \cdot x > 0$ and so $\tan \theta(v_1, r)(q \cdot x) < -v_1 \cdot x$. We obtain the right-hand set in a similar way. Thus, we have that the disagreement region $(v_1 \cdot x)(r \cdot x) < 0$ is a subset of the region $\{|v_1 \cdot x| \le |q \cdot x| \tan \theta(v_1, r)\}$. Since $\tan \theta(v_1, r) \le \theta$ and we have that θ is sufficiently small we can also replace the above region by the larger region: $\{|v_1 \cdot x| \le 2\theta |q \cdot x|\}$. Therefore, we have

$$\begin{split} & \underset{\boldsymbol{x} \sim \mathcal{N}_d}{\mathbf{E}} \left[G^+(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \ \mathbf{1}\{(\boldsymbol{v}_1 \cdot \boldsymbol{x})(\boldsymbol{r} \cdot \boldsymbol{x}) < 0\}\} \right] \\ & \qquad \leq \underset{\boldsymbol{x} \sim \mathcal{N}_d}{\mathbf{E}} \left[G^+(\boldsymbol{v}_2 \cdot \boldsymbol{x}, \dots, \boldsymbol{v}_k \cdot \boldsymbol{x}) \ \mathbf{1}\{|\boldsymbol{v}_1 \cdot \boldsymbol{x}| \leq 2\theta | \boldsymbol{q} \cdot \boldsymbol{x}|\} \right] \,. \end{split}$$

¹⁰⁴⁹ From this point, the proof goes as in the top-1 case. In total, we will get that

$$\begin{split} & \Pr_{\boldsymbol{x}\sim\mathcal{N}_{d}}[T(\boldsymbol{v}_{1}\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_{k}\cdot\boldsymbol{x})\neq T(\boldsymbol{r}\cdot\boldsymbol{x},\boldsymbol{v}_{2}\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_{k}\cdot\boldsymbol{x})] \\ &= \mathop{\mathbf{E}}_{\boldsymbol{x}\sim\mathcal{N}_{d}}\left[\left(G^{+}(\boldsymbol{v}_{2}\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_{k}\cdot\boldsymbol{x})+G^{-}(\boldsymbol{v}_{2}\cdot\boldsymbol{x},\ldots,\boldsymbol{v}_{k}\cdot\boldsymbol{x})\right)|\boldsymbol{q}\cdot\boldsymbol{x}|\,\delta(|\boldsymbol{v}_{1}\cdot\boldsymbol{x}|)\right] \\ &\leq 2\int_{\boldsymbol{x}\in F_{1}}\phi_{d}(\boldsymbol{x})|\boldsymbol{q}\cdot\boldsymbol{x}|d\mu(\boldsymbol{x}) \\ &\leq 2\int_{\boldsymbol{x}\in F_{1}}\phi_{d}(\boldsymbol{x})|\boldsymbol{q}\cdot\boldsymbol{x}|\mathbf{1}\{|\boldsymbol{q}\cdot\boldsymbol{x}|\leq\xi\}d\mu(\boldsymbol{x})+2\int_{\boldsymbol{x}\in F_{1}}\phi_{d}(\boldsymbol{x})|\boldsymbol{q}\cdot\boldsymbol{x}|\mathbf{1}\{|\boldsymbol{q}\cdot\boldsymbol{x}|\geq\xi\}d\mu(\boldsymbol{x}) \\ &\leq 2\xi\int_{\boldsymbol{x}\in F_{1}}\phi_{d}(\boldsymbol{x})d\mu(\boldsymbol{x})+2\int_{\boldsymbol{x}\in F_{1}}\phi_{d}(\boldsymbol{x})|\boldsymbol{q}\cdot\boldsymbol{x}|\mathbf{1}\{|\boldsymbol{q}\cdot\boldsymbol{x}|\geq\xi\}d\mu(\boldsymbol{x}), \end{split}$$

where $d\mu(\mathbf{x})$ is the standard surface measure in \mathbb{R}^d . Let us explain the first inequality above. Note that the space induced by $G^-(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})$ contains the space induced by $G^+(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})$. Hence, in the integration, we can integrate over the surface $F_1 = T(\mathbf{V}\mathbf{x}) \cap \mathbf{1}\{\mathbf{x} : \mathbf{v}_1 \cdot \mathbf{x} = 0\}$ twice. Essentially, this surface corresponds to $\mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} = 0\} \cdot \sum_{S \subseteq [k] \setminus \{1\} : |S| \le r-1} \prod_{i \in S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \ge 0\}$ $0\} \prod_{i \notin S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \le 0\}$. Applying the steps of the top-1 case, we can obtain the desired bound in terms of the Gaussian surface area of F_1 .

Next, for fixed $j \in [k]$, we can apply the above claim sequentially (as we did in the end of the top-1 case) to get

$$\lim_{\theta \to 0} \frac{\mathbf{Pr}_{\boldsymbol{x} \sim \mathcal{N}_d}[T_j(\boldsymbol{U}\boldsymbol{x}) \neq T_j(\boldsymbol{V}\boldsymbol{x})]}{\theta} \le c \cdot \sum_{i \in [k]} \Gamma(F_i^j) \sqrt{\log\left(\frac{1}{\Gamma(F_i^j)} + 1\right)},$$

1058 for some small constant c > 0.

1059 C.2 The proof of Lemma 17

Using the above result, we get that it suffices to control the value $\Gamma(F_i^j)$, where F_i^j is the surface of $T_j(\mathbf{V}\mathbf{x}) \cap \{\mathbf{x} : \mathbf{V}_i \cdot \mathbf{x} = \mathbf{V}_j \cdot \mathbf{x}\}$ for the matrix \mathbf{V} and $i, j \in [k]$. We next have to control the Gaussian surface area of the induced shape, i.e., the quantity

$$\Gamma(\{\boldsymbol{x}: j \in \sigma_{1..r}(\boldsymbol{V}\boldsymbol{x})\} \cap \{\boldsymbol{x}: \boldsymbol{V}_i \cdot \boldsymbol{x} = \boldsymbol{V}_j \cdot \boldsymbol{x}\}).$$

1063 To this end, we give the next lemma.

Lemma 18. Let $r \leq k$ with $r, k \in \mathbb{N}$. For any matrix $V \in \mathbb{R}^{k \times d}$ and $i, j \in [k]$, there exists a matrix $Q = Q^{(i)} \in \mathbb{R}^{k \times d}$ which depends only on i such that

$$\Gamma(F_i^j) := \Gamma(\{\boldsymbol{x} : j \in \sigma_{1..r}(\boldsymbol{V}\boldsymbol{x})\} \cap \{\boldsymbol{x} : \boldsymbol{V}_i \cdot \boldsymbol{x} = \boldsymbol{V}_j \cdot \boldsymbol{x}\}) \leq 2 \cdot \Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[j \in \sigma_{1..r}(\boldsymbol{Q}\boldsymbol{x})].$$

Before proving this result, let us see how to apply it in order to get Lemma 17. We will have that

$$\begin{split} \sum_{i \in [k]} \sum_{j \in [k]} \Gamma(F_i^j) &= \sum_{i \in [k]} \sum_{j \in [k]} \Gamma(\{\boldsymbol{x} : j \in \sigma_{1..r}(\boldsymbol{V}\boldsymbol{x})\} \cap \{\boldsymbol{x} : \boldsymbol{V}_i \cdot \boldsymbol{x} = \boldsymbol{V}_j \cdot \boldsymbol{x}\}) \\ &\leq 2 \sum_{i \in [k]} \sum_{j \in [k]} \Pr_{\boldsymbol{x} \sim \mathcal{N}_d} [j \in \sigma_{1..r}(\boldsymbol{Q}^{(i)}\boldsymbol{x})] \\ &= 2 \sum_{i \in [k]} \sum_{\boldsymbol{x} \sim \mathcal{N}_d} [|\sigma_{1..r}(\boldsymbol{Q}^{(i)}\boldsymbol{x})|] \\ &= 2 \sum_{i \in [k]} r \\ &= 2kr \,. \end{split}$$

Proof of Lemma 18. For this proof, we fix $i, j \in [k]$. The first step is to design the matrix Q. As a first observation, we can subtract the vector V_i from each weight vector and do not affect the resulting orderings. Second, we can assume that the weight vectors that correspond to indices which j beats are unit. Let us be more specific Assume that initially we have that

$$(V_i - V_\ell) \cdot x \ge 0$$

1071 The first observation gives that

$$(V_i - V_i) \cdot \boldsymbol{x} \geq (V_\ell - V_i) \cdot \boldsymbol{x}$$
.

Let us set \widetilde{Q} the intermediate matrix with rows $V_j - V_i$. The second observation states that the inequalities where j beats some index ℓ are not affected by normalization. Note that $\widetilde{Q}_j \cdot x = 0$ and hence $\widetilde{Q}_{\ell} \cdot x \leq 0$. Hence, dividing with non-negative numbers will not affect the order of these two values, i.e.,

$$rac{oldsymbol{Q}_j \cdot oldsymbol{x}}{\|oldsymbol{\widetilde{Q}}_i\|_2} \geq rac{oldsymbol{Q}_\ell \cdot oldsymbol{x}}{\|oldsymbol{\widetilde{Q}}_\ell\|_2}$$

Note that the above ordering is x-dependent, since the indices that j beats depend on x. However, we can normalize any row of \tilde{Q} without affecting the fact that the element j is top-r (since the sign of the inner products is not affected by normalization). This transformation yields a matrix $Q = Q^{(i)}$ and depends only on i (crucially, it is independent of j). For simplicity, we will omit the index i in what follows. For this matrix, we have that

$$\{x: j \in \sigma_{1..r}(Qx), Q_j \cdot x = 0\} = \{x: j \in \sigma_{1..r}(Vx), V_i \cdot x = V_j \cdot x\}.$$

1081 We will now prove that

$$\mathbf{\Pr}_{oldsymbol{x} \sim \mathcal{N}_d}[j \in \sigma_{1..r}(oldsymbol{Q} oldsymbol{x})] \geq rac{\Gamma(F_i^j)}{2}$$
 .

< _ i .

Let us fix some x and set $x^{\parallel} = \operatorname{proj}_{Q_j} x$ and $x^{\perp} = \operatorname{proj}_{Q_j^{\perp}} x$. We assume that x lies in the set $x = x = \frac{1}{2} \{x : j \in \sigma_{1..r}(Qx)\}$. This implies that there exist an index set I of size at least k - r so that if $\ell \in I$ then

$$oldsymbol{Q}_j \cdot oldsymbol{x}^{\parallel} + oldsymbol{Q}_j \cdot oldsymbol{x}^{\perp} \geq oldsymbol{Q}_\ell \cdot oldsymbol{x}^{\parallel} + oldsymbol{Q}_\ell \cdot oldsymbol{x}^{\perp}$$
 .

1085 Let us condition on the event

$$oldsymbol{Q}_j \cdot oldsymbol{x}^\perp \geq oldsymbol{Q}_\ell \cdot oldsymbol{x}^\perp$$
 ,

1086 We hence get that

$$oldsymbol{Q}_j \cdot oldsymbol{x}^{\parallel} = (oldsymbol{Q}_j \cdot oldsymbol{Q}_j) \cdot (oldsymbol{Q}_j \cdot oldsymbol{x}) \geq oldsymbol{Q}_\ell \cdot oldsymbol{x}^{\parallel} = (oldsymbol{Q}_\ell \cdot oldsymbol{Q}_j) \cdot (oldsymbol{Q}_j \cdot oldsymbol{x})$$

Using that Q_j is unit, that the inner product between Q_ℓ and Q_j is at most one and that $Q_j \cdot x$ is a univariate Gaussian, we get that

$$\Pr_{z \sim \mathcal{N}(0,1)} [z \cdot (1 - \boldsymbol{Q}_{\ell} \cdot \boldsymbol{Q}_j) \ge 0] = 1/2.$$

1089 The above discussion implies that

$$\Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[j \in \sigma_{1..r}(\boldsymbol{Q}\boldsymbol{x})] = \Pr_{\boldsymbol{x} \sim \mathcal{N}_d}[(\forall \ell \in I) \; \boldsymbol{Q}_j \cdot \boldsymbol{x}^{\parallel} + \boldsymbol{Q}_j \cdot \boldsymbol{x}^{\perp} \geq \boldsymbol{Q}_\ell \cdot \boldsymbol{x}^{\parallel} + \boldsymbol{Q}_\ell \cdot \boldsymbol{x}^{\perp}]$$

1090 and so $\mathbf{Pr}_{m{x}\sim\mathcal{N}_d}[j\in\sigma_{1..r}(m{Q}m{x})]$ equals to

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[(\forall \ell\in I) \, \boldsymbol{Q}_j \cdot \boldsymbol{x}^{\parallel} \geq \boldsymbol{Q}_j \cdot \boldsymbol{x}^{\parallel} \mid (\forall \ell\in I) \, \boldsymbol{Q}_j \cdot \boldsymbol{x}^{\perp} \geq \boldsymbol{Q}_\ell \cdot \boldsymbol{x}^{\perp}] \cdot \Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[(\forall \ell\in I) \, \boldsymbol{Q}_j \cdot \boldsymbol{x}^{\perp} \geq \boldsymbol{Q}_\ell \cdot \boldsymbol{x}^{\perp}].$$

However, in the above product, we have that the first term is 1/2 and the second term is the probability that $j \in \sigma_{1..r}(Qx^{\perp})$, i.e.,

$$\Pr_{\boldsymbol{x}\sim\mathcal{N}_d}[j\in\sigma_{1..r}(\boldsymbol{Q}\boldsymbol{x})]\geq\frac{\Pr[j\in\sigma_{1..r}(\boldsymbol{Q}\boldsymbol{x}^{\perp})]}{2}=\Gamma(F_i^j)/2\,,$$

since the space in the RHS is low-dimensional and corresponds to the desired surface.

1094 D Distribution-Free Lower Bounds for Top-1 Disagreement Error

We begin with some definitions concerning the PAC Label Ranking setting. Let \mathcal{X} be an instance space and $\mathcal{Y} = \mathbb{S}_k$ be the space of labels, which are rankings over k elements. A sorting function or hypothesis is a mapping $h : \mathcal{X} \to \mathbb{S}_k$. We denote by $h_1(x)$ the top-1 element of the ranking h(x). A hypothesis class is a set of classifiers $\mathcal{H} \subset \mathbb{S}_k^{\mathcal{X}}$.

Top-1 Disagreement Error. The top-1 disagreement error with respect to a joint distribution \mathcal{D} over $\mathcal{X} \times \mathbb{S}_k$ equals to the probability $\mathbf{Pr}_{(x,\sigma)\sim\mathcal{D}}[h_1(x) \neq \sigma^{-1}(1)]$. We mainly consider learning in the **realizable** case, which means that there is $h^* \in \mathcal{H}$ which has (almost surely) zero error. Therefore, we can focus on the marginal distribution \mathcal{D}_x over \mathcal{X} and denote the top-1 disagreement error of a sorting function h with respect to the true hypothesis h^* by $\operatorname{Err}_{\mathcal{D}_x,h^*}(h) := \mathbf{Pr}_{x\sim\mathcal{D}_x}[h_1(x) \neq h_1^*(x)]$.

A learning algorithm is a function \mathcal{A} that receives a training set of m instances, $S \in \mathcal{X}^m$, together with their labels according to h^* . We denote the restriction of h^* to the instances in S by $h^*|_S$. The output of the algorithm \mathcal{A} , denoted $\mathcal{A}(S, h^*|_S)$ is a sorting function. A learning algorithm is proper if it always outputs a hypothesis from \mathcal{H} .

The top-1 PAC Label Ranking sample complexity of a learning algorithm \mathcal{A} is the function $m_{\mathcal{A},\mathcal{H}}^{(1)}$ defined as follows: for every $\epsilon, \delta > 0, m_{\mathcal{A},\mathcal{H}}^{(1)}(\epsilon, \delta)$ is the minimal integer such that for every $m \geq m_{\mathcal{A},\mathcal{H}}^{(1)}(\epsilon, \delta)$, every distribution \mathcal{D}_x on \mathcal{X} , and every target hypothesis $h^* \in \mathcal{H}$, **Pr**_{S~ \mathcal{D}_x^m} [Err $_{\mathcal{D}_x,h^*}(\mathcal{A}(S,h^*|_S)) > \epsilon$] $\leq \delta$. In this case, we say that the learning algorithm (ϵ, δ) learns the class of sorting functions \mathcal{H} with respect to the top-1 disagreement error. If no integer satisfies the inequality above, define $m_{\mathcal{A}}^{(1)}(\epsilon, \delta) = \infty$. \mathcal{H} is learnable with \mathcal{A} if for all ϵ and δ the sample complexity is finite. The **top-1 PAC Label Ranking sample complexity** of a class \mathcal{H} is $m_{\text{PAC},\mathcal{H}}^{(1)}(\epsilon, \delta) = \inf_{\mathcal{A}} m_{\mathcal{A},\mathcal{H}}^{(1)}(\epsilon, \delta)$, where the infimum is taken over all learning algorithms. Clearly, the above top-1 definition can be extended to the top-r setting.

In this section, we show the next result. We denote by $\mathcal{L}_{d,k}$ the class of Linear Sorting functions in *d* dimensions with *k* labels.

Theorem 7. In the realizable PAC Label Ranking setting, any algorithm that (ϵ, δ) -learns the class $\mathcal{L}_{d,k}$ with respect to the top-1 disagreement error requires at least $\Omega((dk + \log(1/\delta))/\epsilon)$ samples.

1121 D.1 Top-1 Ranking Natarajan Dimension

In order to establish the above result, we introduce a variant of the standard Natarajan dimension [Nat89, BDCBL92, DSBDSS11, DSS14]. For a ranking π , we will also let $L_1(\pi)$ its top-1 element and $L_{3..k}(\pi)$ the ranking after deleting its top-2 part.

Definition 3 (Top-1 Ranking Natarajan Dimension). Let $\mathcal{H} \subseteq \mathbb{S}_k^{\mathcal{X}}$ be a hypothesis class of sorting 1125 functions and let $S \subseteq \mathcal{X}$. We say that \mathcal{H} N-shatters S if there exist two mappings $f_1, f_2: S \to \mathbb{S}_k$ 1126 such that for every $y \in S$, $L_1(f_1(y)) \neq L_1(f_2(y))$ and $L_{3..k}(f_1(y)) = L_{3..k}(f_2(y))$ and for every 1127 1128

 $T \subseteq S$, there exists a sorting function $g \in \mathcal{H}$ such that

(i) $\forall x \in T$, $g(x) = f_1(x)$, and (ii) $\forall x \in S \setminus T$, $g(x) = f_2(x)$.

The top-1 Ranking Natarajan dimension of \mathcal{H} , denoted $d_N^{(1)}(\mathcal{H})$ is the maximal cardinality of a set 1129 that is N-shattered by \mathcal{H} . 1130

First, we connect LR PAC learnability to the top-1 disagreement error with the notion of top-1 ranking 1131 Natarajan dimension. 1132

Theorem 8 (Top-1-Natarajan Lower Bounds Sample Complexity). In the realizable PAC Label 1133 *Ranking setting, we have for every hypothesis class* $\mathcal{H} \subseteq \mathbb{S}_{k}^{\mathcal{X}}$ 1134

$$m_{\text{PAC},\mathcal{H}}^{(1)}(\epsilon,\delta) = \Omega\left(\frac{d_N^{(1)}(\mathcal{H}) + \ln(1/\delta)}{\epsilon}\right).$$

Proof. Let $\mathcal{H} \subseteq \mathbb{S}_k^{\mathcal{X}}$ be a hypothesis of sorting functions of top-1-Natarajan dimension $d_N^{(1)} = d_N$. Consider the binary hypothesis class $\mathcal{H}_{\text{bin}} = \{0,1\}^{[d_N]}$ which contains all the classifiers from 1135 1136 $[d_N] = \{1, ..., d_N\}$ to $\{0, 1\}$. It suffices to show the following. 1137

Claim 12. It holds that $m_{\text{PAC},\mathcal{H}}^{(1)}(\epsilon,\delta) \geq m_{\text{PAC},\mathcal{H}_{\text{bin}}}(\epsilon,\delta)$. 1138

This is sufficient since we have that $m_{\text{PAC},\mathcal{H}_{\text{bin}}}(\epsilon,\delta) = \Omega\left(\frac{\operatorname{VC}(\mathcal{H}_{\text{bin}}) + \ln(1/\delta)}{\epsilon}\right)$ and $\operatorname{VC}(\mathcal{H}_{\text{bin}}) = d_N$. 1139 Let us now prove the claim. 1140

We assume that the instance space is the set \mathcal{X} . Assume that A is a learning algorithm for the 1141 hypothesis class $\mathcal{H} \subseteq \mathbb{S}_k^{\mathcal{X}}$ and A_{bin} is a learning algorithm for the associated binary class \mathcal{H}_{bin} . It 1142 suffices to show that A requires at least as many samples as $A_{\rm bin}$. In fact, we will show that whenever 1143 A_{bin} errs, so does A. Let $S = \{s_1, ..., s_{d_N}\}, f_0, f_1$ be the set and the two functions that witness that the top-1-Natarajan dimension of \mathcal{H} is d_N . Given a training set $(x_i, y_i)_{i \in [m]} \in ([d_N] \times \{0, 1\})^m$, we 1144 1145 set $g: \mathcal{X} \to \mathbb{S}_k$ be equal to the output of the algorithm A with input $(s_{x_i}, f_{y_i}(x_i))_{i \in [m]} \in (S \times \mathbb{S}_k)^m$. 1146 We also set f be the output of the algorithm A_{bin} with input $(x_i, y_i)_{i \in [m]}$ by setting f(i) = 11147 if and only if $L_1(g(s_i)) = L_1(f_1(s_i))$. We will show that whenever A_{bin} errs, so does A. Fix 1148 $(x_i, y_i) \in S \times \{0, 1\}$. Assume that $A_{\text{bin}}(x_i) \neq y_i$ and say $y_i = 0$. Then f(i) = 1 and so $L_1(g(s_i)) = L_1(f_1(s_i)) \neq L_1(f_0(s_i))$. This implies that A errs. The case $y_i = 1$ is similar. \Box 1149 1150

D.2 Lower Bound for top-1 disagreement error for LSFs 1151

Theorem 9 (Top-1 Natarajan Dimension of LSFs). Consider the hypothesis class $\mathcal{L}_{d,k} = \{\sigma_W :$ 1152 $\mathbb{R}^d \to \mathbb{S}_k : \sigma_{\boldsymbol{W}}(\boldsymbol{x}) = \operatorname{argsort}(\boldsymbol{W}\boldsymbol{x}), \boldsymbol{W} \in \mathbb{R}^{k \times d} \}.$ Then, $d_N^{(1)}(\mathcal{L}_{d,k}) = \Omega(dk).$ 1153

Proof. Fix $k \in \mathbb{N}$. Let us consider the case d = 2 that will correspond as the building block for 1154 the general case d > 2. Let us first choose the set of points: Set P be the collection of pairs 1155 $P = \{(2i-1,2i)\}_{i\in[b]}$ for any $i \in [b]$ with $b = \lfloor k/2 \rfloor$ and $S = \{x_m\}_{m \in P}$ where these points 1156 correspond to |P| equidistributed points on the unit sphere in \mathbb{R}^2 . This set of points has size 1157 $|P| = \Theta(k)$ and we are going to N-shatter it using $\mathcal{L}_{2,k}$. 1158

Consider the matrix $W \in \mathbb{R}^{k \times 2}$ so that $\{W_i\}_{i \in [k]}$ correspond to the rows of W. The structure of 1159 the problem relies on the hyperplanes with normal vectors $(W_i - W_j)_{i \neq j}$ and our choice of W will 1160 rely on these hyperplanes. For any m = (2i - 1, 2i), we set W_{2i-1}, W_{2i} on the unit sphere so that 1161 $W_{2i-1} \cdot W_{2i} = 1 - \phi$ with $\phi \in (0,1)$ sufficiently small (set $\arccos(1-\phi) = 2\pi/(100k)$) and let 1162 C_m be the cone generated by these two vectors with axis I_m . We place W_{2i-1} so that the distance 1163 between x_m and the hyperplane I_m is sufficiently small (say that the angle between x_m and I_m is 1164 $\arccos(1-\phi)/100$). Note that the normal vector of I_m is $W_{2i-1} - W_{2i}$ and we place x_m so that 1165 it has positive correlation with this vector. This uniquely identifies the location of W_{2i} . Crucially, 1166 each vector x_m has the following properties: (i) x_m is very close to the boundary of the hyperplane 1167

with normal vector $(W_{2i-1} - W_{2i})$, (ii) $W_{2i-1} \cdot x_m > W_{2i} \cdot x > W_j \cdot x_m$ for any $j \notin m$ and 1168 (iii) x_m is far from any boundary induced by hyperplanes with normal vectors $W_j - W_{j'}$ for any 1169 $(j, j') \neq m.$ 1170

Since the points are well-separated on the unit sphere, for any $m = (2i - 1, 2i) \in P$, we have 1171 $W_{2i-1} \cdot W_{2i} = 1 - \phi \approx 1$ and for any other pair of indices $(i, j) \notin P$, there exists $c = c(k) \in (0, 1)$, 1172 $|\langle \boldsymbol{W}_i, \boldsymbol{W}_j \rangle| \leq c.$ 1173

For any $m = (2i - 1, 2i) \in P$, we set $W'_{2i-1} - W'_{2i} = R_{\theta}(W_{2i-1} - W_{2i})$ for some θ to be chosen, where R_{θ} is the 2 × 2 rotation matrix. We choose θ so that each point x_m for $m = (2i - 1, 2i) \in P$ 1174 1175 with $(W_{2i-1} - W_{2i}) \cdot x_m > 0$ satisfies $(W'_{2i-1} - W'_{2i}) \cdot x_m < 0$. The main idea is that since x_m 1176 has the properties (i)-(iii) described above, the rankings induced by the vectors $W x_m$ and $W' x_m$ 1177 will be different in the first two positions but the same in the rest. 1178

Given the training set $\{x_m\}_{m \in P}$, we have to construct f_0, f_1 and verify that they satisfy the top-1 Ranking Natarajan conditions. For m = (2i - 1, 2i), we have that $f_0(x_m) = (2i - 1, 2i, \pi)$ 1179 1180 and $f_1(\boldsymbol{x}_m) = (2i, 2i - 1, \pi)$ for some ranking π of size k - 2 that depends on m. Specifically, 1181 we will set $f_0(x) = \sigma(Wx)$ and $f_1(x) = \sigma(W'x)$, where σ gives the decreasing ordering of 1182 the elements of the input vector. By the choice of the set S and W, W', it remains to show 1183 that the k-2 last elements of the rankings $f_0(x_m)$ (say π_0) and of $f_1(x_m)$ (say π_1) are in the 1184 same order, i.e., $L_{3..k}(f_0(\boldsymbol{x}_m)) = L_{3..k}(f_1(\boldsymbol{x}_m))$. Assume that $u \succ v$ in π_0 . It suffices to show 1185 that $(W'_u - W'_v) \cdot x_m \ge 0$, i.e., the order of u and v is preserved when transforming W to 1186 W'. We have that $(W_u - W_v) \cdot x_m > c_1$ for some constant $c_1 > 0$ (c_1 is the minimum over 1187 $(u, v) \neq m = (2i - 1, 2i)$. Hence, we can pick θ small enough so that $(\mathbf{W}'_u - \mathbf{W}'_v) \cdot \mathbf{x}_m > c_2$ and 1188 this can be done for any pair u, v that does not correspond to m. This implies that $\pi_0 = \pi_1 = \pi$. In 1189 particular, we have that 1190

$$(\boldsymbol{W}'_{u} - \boldsymbol{W}'_{v}) \cdot \boldsymbol{x}_{m} = \cos(\theta) \cdot (\boldsymbol{W}_{u} - \boldsymbol{W}_{v}) \cdot \boldsymbol{x}_{m} + \sin(\theta) \cdot (W^{(1)}_{uv} x^{(2)}_{m} - W^{(2)}_{uv} x^{(1)}_{m}) > c_{2} > 0$$

for some θ sufficiently small, where $W_{uv}^{(t)}$ is the *t*-th entry of the vector $W_u - W_v$ for $t \in \{1, 2\}$ and 1191 $\boldsymbol{x}_m, \boldsymbol{W}_u, \boldsymbol{W}_v$ are unit vectors. 1192

For any subset T of S, it remains to choose a linear classifier in $\mathcal{L}_{2,k}$ (which is allowed to depend 1193 on T). For any $T \subseteq S = \{x_m\}_{m \in P}$, we consider the matrix $\overline{W} \in \mathbb{R}^{k \times 2}$ so that for the *i*-th row 1194 $\overline{W}_i = W_i \{i \in m \in T\} + W'_i \{i \in m \in S \setminus T\}$ for any $i \in [k]$. This is valid since the pairs 1195 $m \in P$ partition [k]. We have to show the following two properties: (i) $\sigma(\overline{W}x) = f_0(x)$ for $x \in T$ 1196 and (ii) $\sigma(\overline{W}x) = f_1(x)$ for $x \in S \setminus T$. 1197

Assume that m = (2i - 1, 2i) and $\boldsymbol{x}_m \in T$. We have that $f_0(\boldsymbol{x}_m) = (2i - 1, 2i, \pi)$ and $\overline{\boldsymbol{W}}_{2i-1}$ 1198 $\overline{W}_{2i} = W_{2i-1} - W_{2i}$ and so $2i - 1 \succ 2i$ in the ranking $\sigma(\overline{W}x_m)$. It remains to show that the 1199 remaining $\binom{k}{2} - 1$ pairwise comparisons are the same in the two rankings. Let us consider a pair of 1200 points $u \neq v$ so that $u \succ v$ in $f_0(\mathbf{x}_m)$. It suffices to show that $u \succ v$ in $\sigma(\overline{\mathbf{W}}\mathbf{x}_m)$. 1201

1. If u, v are so that
$$\overline{W}_u - \overline{W}_v = W_u - W_v$$
, the result holds.

1203

2. If u, v are so that $\overline{W}_u - \overline{W}_v = W_u - W'_v$: In this case, u and v lie in a different pair of P and this implies that the correct direction is preserved if θ is appropriately chosen. For 1204 θ as above, it holds that $(W_u - R_{\theta}W_v) \cdot x_m$ has the same sign as $(W_u - W_v) \cdot x_m$. In 1205 particular, 1206

$$W_{u} \cdot x_{m} - R_{\theta} W_{v} \cdot x_{m} = W_{u} \cdot x_{m} - (\cos(\theta) W_{v}^{(1)} - \sin(\theta) W_{v}^{(2)}) x_{m}^{(1)} - (\sin(\theta) W_{v}^{(1)} + \cos(\theta) W_{v}^{(2)}) x_{m}^{(2)}$$

1207

and so

$$(\boldsymbol{W}_{u} - \boldsymbol{W}_{v}') \cdot \boldsymbol{x}_{m} = \cos(\theta) \cdot (\boldsymbol{W}_{u} - \boldsymbol{W}_{v}) \cdot \boldsymbol{x}_{m} + \sin(\theta)(W_{v}^{(2)}x_{m}^{(1)} - W_{v}^{(1)}x_{m}^{(2)}) > 0.$$

3. If u, v are so that $\overline{W}_u - \overline{W}_v = W'_u - W'_v$, the analysis for the inner product with x_m will 1208 be similar. 1209

We now have to extend this proof for d > 2. We will "tensorize" the above construction as follows. 1210 Let $S = \{y_{m_i}\}_{m \in [b], i \in [d/2]}$ with $|S| = \lfloor k/2 \rfloor \cdot \lfloor d/2 \rfloor$. We first define the points of S: For $s \in [d]$, 1211

set $y_{mj}[s] = x_m[1]1\{s = 2j - 1\} + x_m[2]1\{s = 2j\}$ with $y_{mj} \in \mathbb{R}^d$, i.e., y_{mj} has the values of x_m at the consecutive entries indicated by $m = (2i - 1, 2i) \in P$ and zeros at the other positions.

We have to show that the set S is N-shattered. Given $T \subseteq S$, we are going to create the matrix $\overline{W} \in \mathbb{R}^{k \times d}$. For illustration, think of each row of the matrix as having d/2 blocks of size two. If $y_{mj} \in T$ with m = (2i-1, 2i), set the two associated rows (indicated by m) of \overline{W} with W_{2i-1}, W_{2i} at the j-th block and with W'_{2i-1}, W'_{2i} otherwise. We will have that $\sigma(\overline{W}y) = f_0(y)$ if $y \in T$ and $\sigma(\overline{W}y) = f_1(y)$ otherwise and the analysis is the same as the d = 2 case.

1219 E Examples of Noisy Ranking Distributions

Definition 4 (Mallows model [Mal57]). Consider k alternatives and let $\pi \in S_k, \phi \in [0, 1]$. The Mallows distribution $\mathcal{M}_{Mal}(\pi, \phi)$ with central ranking π and spread parameter ϕ is a probability measure over S_k with density $\mathbf{Pr}_{\sigma \sim \mathcal{M}_{Mal}(\pi, \phi)}[\sigma]$ that is proportional to $\phi^{d(\sigma, \pi)}$, where d is a ranking distance.

We focus on Mallows models accociated with the Kendall's Tau distance $d = d_{KT}$ (the standard distance, not the normalized one), which measures the number of discordant pairs.

Fact 2. When $\phi < 1$, the Mallows model $\mathcal{M}_{Mal}(\pi, \phi)$ is a ranking distribution with bounded noise at most $\frac{1+\phi}{4} < 1/2$.

1228 Proof. The following property holds [Mal57]

$$\Pr_{\sigma \sim \mathcal{M}_{\text{Mal}}(\pi,\phi)}[\sigma(i) < \sigma(j) | \pi(i) < \pi(j)] = \frac{\pi(j) - \pi(i) + 1}{1 - \phi^{\pi(j) - \pi(i) + 1}} - \frac{\pi(j) - \pi(i)}{1 - \phi^{\pi(j) - \pi(i)}} \ge \frac{1}{2} + \frac{1 - \phi}{4}.$$

1229

The Bradley-Terry-Luce model [BT52, Luc12] is the most studied pairwise comparisons model. In his seminal paper, Mallows [Mal57] also studied the following natural ranking distribution:

1232 **Definition 5** (Bradley-Terry-Mallows [Mal57]). Consider a score vector $w \in \mathbb{R}^k_+$ with k distinct 1233 entries and let π be the ranking induced by the values of w in decreasing order. The Bradley-Terry-1234 Mallows distribution $\mathcal{M}_{BTM}(w)$ with central ranking π is a probability measure over \mathbb{S}_k with density 1235 $\mathbf{Pr}_{\sigma \sim \mathcal{M}_{BTM}(w)}[\sigma]$ that is proportional to $\prod_{i \succ_{\sigma} j} \frac{w_i}{w_i + w_j}$.

Lemma 19. There exists a real number $0 < \eta < 1/2$ so that the Bradley-Terry-Mallows distribution $\mathcal{M}_{BTM}(w)$ is a ranking distribution with bounded noise at most η .

Proof. In the standard Bradley-Terry-Luce model, the pairwise comparison between the alternatives *i*, *j* is a Bernoulli random variable with $\mathbf{Pr}[i \succ j] = w_i/(w_i + w_j)$. The Bradley-Terry-Mallows distribution can be considered as the Bradley-Terry-Luce model conditioned on the event that all the pairwise comparisons are consistent to a ranking. Hence, we have that

$$\Pr_{\sigma \sim \mathcal{M}_{\text{BTM}}(\boldsymbol{w})}[\sigma] = \frac{1}{Z(k, \boldsymbol{w})} \prod_{i \succ_{\sigma j}} \frac{w_i}{w_i + w_j} \,.$$

Let us set $\mathcal{A}_{i \succ j} = \{ \sigma \in \mathbb{S}_k : \sigma(i) < \sigma(j) \}$. We are interested in the following probability

$$\Pr_{\sigma \sim \mathcal{M}_{\mathrm{BTM}}(\boldsymbol{w})}[i \succ_{\sigma} j | w_i > w_j] = \Pr_{\sigma \sim \mathcal{M}_{\mathrm{BTM}}(\boldsymbol{w})}[\sigma(i) < \sigma(j) | w_i > w_j] = \frac{1}{Z(k, \boldsymbol{w})} \sum_{\sigma \in \mathcal{A}_{i \succ j}} \prod_{p \succ_{\sigma} q} \frac{w_p}{w_p + w_q}$$

1243 Note that in order to show the desired property, it suffices to show that

$$\sum_{\sigma \in \mathcal{A}_{i \succ j}} \prod_{p \succ \sigma q} \frac{w_p}{w_p + w_q} > \sum_{\sigma \in \mathcal{A}_{i \prec j}} \prod_{p \succ \sigma q} \frac{w_p}{w_p + w_q}$$

First, observe that there exists a correspondence mapping $\sigma \in A_{i \succ j}$ to $A_{i \prec j}$, where one flips the elements *i* and *j*. Hence, it suffices to show that the mass of the ranking $(u_a)i(u_b)j(u_c)$ is larger than the one of the ranking $(u_a)j(u_b)i(u_c)$, where u_a, u_b, u_c are permutations of length between 0 and k-2 with elements in $[k] \setminus \{i, j\}$. For the two above rankings, the only terms of the product that are not identical are the following

$$\frac{w_i}{w_i + w_j} \prod_{x \in u_b} \frac{w_i}{w_i + w_x} \frac{w_x}{w_x + w_j} > \frac{w_j}{w_i + w_j} \prod_{x \in u_b} \frac{w_j}{w_j + w_x} \frac{w_x}{w_x + w_i} \,,$$

1249 since $w_i > w_j$ and so the result follows.